

Postgres Pro Shardman 17.5.1 Documentation



Postgres Professional

<https://postgrespro.com>

Postgres Pro Shardman 17.5.1 Documentation

Postgres Professional

Copyright © 2016–2025 The Postgres Professional company

Legal Notice

This documentation is intended solely for the use with the Postgres Pro DBMS and for users of this DBMS.

It is not allowed to use the documentation for third-party products or as part of documentation for other products.

Other terms of use of the documentation are given in the User Agreement.

Postgres Pro is Copyright © 2016–2025 by Postgres Professional.

IN NO EVENT SHALL THE POSTGRES PROFESSIONAL COMPANY BE LIABLE TO ANY PARTY FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, OR CONSEQUENTIAL DAMAGES, INCLUDING LOST PROFITS, ARISING OUT OF THE USE OF POSTGRES PRO DBMS IN ALL VERSIONS AND ITS DOCUMENTATION, EVEN IF THE POSTGRES PROFESSIONAL COMPANY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

THE POSTGRES PROFESSIONAL COMPANY SPECIFICALLY DISCLAIMS ANY WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE POSTGRES PRO DBMS IN ALL VERSIONS AND ITS DOCUMENTATION PROVIDED HEREUNDER IS ON AN “AS-IS” BASIS, AND THE POSTGRES PROFESSIONAL COMPANY HAS NO OBLIGATIONS TO PROVIDE MAINTENANCE, SUPPORT, UPDATES, ENHANCEMENTS, OR MODIFICATIONS.

Distribution of this documentation or its parts that are not contained in the PostgreSQL documentation, in the original or modified form, requires an explicit written permission from the Postgres Professional company.

Postgres Pro DBMS documentation is based on the PostgreSQL documentation, which is distributed under the following license:

PostgreSQL is Copyright © 1996–2025 by the PostgreSQL Global Development Group.

Postgres95 is Copyright © 1994–5 by the Regents of the University of California.

Permission to use, copy, modify, and distribute this software and its documentation for any purpose, without fee, and without a written agreement is hereby granted, provided that the above copyright notice and this paragraph and the following two paragraphs appear in all copies.

IN NO EVENT SHALL THE UNIVERSITY OF CALIFORNIA BE LIABLE TO ANY PARTY FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, OR CONSEQUENTIAL DAMAGES, INCLUDING LOST PROFITS, ARISING OUT OF THE USE OF THIS SOFTWARE AND ITS DOCUMENTATION, EVEN IF THE UNIVERSITY OF CALIFORNIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

THE UNIVERSITY OF CALIFORNIA SPECIFICALLY DISCLAIMS ANY WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE SOFTWARE PROVIDED HEREUNDER IS ON AN “AS-IS” BASIS, AND THE UNIVERSITY OF CALIFORNIA HAS NO OBLIGATIONS TO PROVIDE MAINTENANCE, SUPPORT, UPDATES, ENHANCEMENTS, OR MODIFICATIONS.

Preface	xxii
1. What Is Postgres Pro Shardman?	xxii
2. Difference between Postgres Pro Shardman and PostgreSQL	xxiv
3. A Brief History of PostgreSQL	xxv
3.1. The Berkeley POSTGRES Project	xxv
3.2. Postgres95	xxvi
3.3. PostgreSQL	xxvi
4. Conventions	xxvii
5. Bug Reporting Guidelines	xxvii
5.1. Identifying Bugs	xxvii
5.2. What to Report	xxvii
5.3. Where to Report Bugs	xxviii
I. Tutorial	1
1. Getting Started	2
1.1. Quickstart Guide	2
1.2. When to use	10
1.3. Architectural Fundamentals	11
1.4. Creating a Database	11
1.5. Accessing a Database	13
2. The SQL Language	16
2.1. Introduction	16
2.2. Concepts	16
2.3. Creating a New Table	16
2.4. Populating a Table With Rows	18
2.5. Querying a Table	18
2.6. Joins Between Tables	20
2.7. Aggregate Functions	21
2.8. Updates	23
2.9. Deletions	23
3. Advanced Features	25
3.1. Introduction	25
3.2. Views	25
3.3. Foreign Keys	25
3.4. Transactions	26
3.5. Window Functions	27
3.6. Inheritance	30
3.7. Conclusion	31
II. The SQL Language	32
4. SQL Syntax	33
4.1. Lexical Structure	33
4.2. Value Expressions	41
4.3. Calling Functions	53
5. Data Definition	56
5.1. Table Basics	56
5.2. Default Values	57
5.3. Identity Columns	58
5.4. Generated Columns	59
5.5. Constraints	60
5.6. System Columns	68
5.7. Modifying Tables	69
5.8. Privileges	71
5.9. Row Security Policies	75
5.10. Non-Global Schemas	81
5.11. Inheritance	86
5.12. Table Partitioning	89
5.13. Foreign Data	101
5.14. Other Database Objects	102
5.15. Dependency Tracking	102

6. Data Manipulation	104
6.1. Inserting Data	104
6.2. Updating Data	105
6.3. Deleting Data	105
6.4. Returning Data from Modified Rows	106
7. Queries	107
7.1. Overview	107
7.2. Table Expressions	107
7.3. Select Lists	121
7.4. Combining Queries (UNION, INTERSECT, EXCEPT)	122
7.5. Sorting Rows (ORDER BY)	123
7.6. LIMIT and OFFSET	124
7.7. VALUES Lists	124
7.8. WITH Queries (Common Table Expressions)	125
8. Data Types	133
8.1. Numeric Types	134
8.2. Monetary Types	139
8.3. Character Types	140
8.4. Binary Data Types	142
8.5. Date/Time Types	144
8.6. Boolean Type	153
8.7. Enumerated Types	154
8.8. Geometric Types	155
8.9. Network Address Types	157
8.10. Bit String Types	160
8.11. Text Search Types	161
8.12. UUID Type	163
8.13. XML Type	163
8.14. JSON Types	165
8.15. Arrays	174
8.16. Composite Types	182
8.17. Range Types	188
8.18. Domain Types	192
8.19. Object Identifier Types	193
8.20. pg_lsn Type	195
8.21. Pseudo-Types	196
9. Functions and Operators	198
9.1. Logical Operators	198
9.2. Comparison Functions and Operators	198
9.3. Mathematical Functions and Operators	202
9.4. String Functions and Operators	210
9.5. Binary String Functions and Operators	219
9.6. Bit String Functions and Operators	223
9.7. Pattern Matching	224
9.8. Data Type Formatting Functions	242
9.9. Date/Time Functions and Operators	249
9.10. Enum Support Functions	264
9.11. Geometric Functions and Operators	264
9.12. Network Address Functions and Operators	271
9.13. Text Search Functions and Operators	274
9.14. UUID Functions	279
9.15. XML Functions	280
9.16. JSON Functions and Operators	292
9.17. Sequence Manipulation Functions	322
9.18. Conditional Expressions	323
9.19. Array Functions and Operators	326
9.20. Range/Multirange Functions and Operators	329

9.21. Aggregate Functions	334
9.22. Window Functions	341
9.23. Merge Support Functions	343
9.24. Subquery Expressions	343
9.25. Row and Array Comparisons	346
9.26. Set Returning Functions	348
9.27. System Information Functions and Operators	352
9.28. System Administration Functions	371
9.29. Trigger Functions	391
9.30. Event Trigger Functions	392
9.31. Statistics Information Functions	395
9.32. Functions for Distributed System	395
10. Type Conversion	399
10.1. Overview	399
10.2. Operators	400
10.3. Functions	403
10.4. Value Storage	407
10.5. UNION, CASE, and Related Constructs	407
10.6. SELECT Output Columns	409
11. Indexes	410
11.1. Introduction	410
11.2. Index Types	411
11.3. Multicolumn Indexes	413
11.4. Indexes and ORDER BY	414
11.5. Combining Multiple Indexes	414
11.6. Unique Indexes	415
11.7. Indexes on Expressions	415
11.8. Partial Indexes	416
11.9. Index-Only Scans and Covering Indexes	418
11.10. Operator Classes and Operator Families	421
11.11. Indexes and Collations	422
11.12. Examining Index Usage	422
11.13. Using k-NN Algorithm for Optimizing Queries	423
12. Full Text Search	426
12.1. Introduction	426
12.2. Tables and Indexes	429
12.3. Controlling Text Search	431
12.4. Additional Features	438
12.5. Parsers	442
12.6. Dictionaries	444
12.7. Configuration Example	452
12.8. Testing and Debugging Text Search	453
12.9. Preferred Index Types for Text Search	457
12.10. psql Support	458
12.11. Limitations	460
13. Concurrency Control	462
13.1. Introduction	462
13.2. Transaction Isolation	462
13.3. Explicit Locking	468
13.4. Data Consistency Checks at the Application Level	474
13.5. Serialization Failure Handling	475
13.6. Caveats	476
13.7. Locking and Indexes	476
14. Performance Tips	477
14.1. Using EXPLAIN	477
14.2. Statistics Used by the Planner	489
14.3. Controlling the Planner with Explicit JOIN Clauses	494

14.4. Populating a Database	495
14.5. Autoprepared Statements	497
14.6. Postgres Pro Shardman Performance Tuning	498
15. Parallel Query	510
15.1. How Parallel Query Works	510
15.2. When Can Parallel Query Be Used?	511
15.3. Parallel Plans	511
15.4. Parallel Safety	513
16. Autonomous Transactions	515
16.1. Behavior	515
16.2. Visibility	516
16.3. SQL Grammar Extension for Autonomous Transactions	517
16.4. PL/pgSQL Grammar Extension for Autonomous Transactions	517
16.5. PL/Python Extension for Autonomous Transactions	518
III. Server Administration	519
17. Binary Installation	520
17.1. Installing Postgres Pro Enterprise	520
17.2. Installing Additional Supplied Modules	528
17.3. Migration	530
18. Server Setup and Operation	558
18.1. The Postgres Pro User Account	558
18.2. Creating a Database Cluster	558
18.3. Starting the Database Server	560
18.4. Managing Kernel Resources	564
18.5. Shutting Down the Server	571
18.6. Upgrading a Postgres Pro Cluster	572
18.7. Preventing Server Spoofing	574
18.8. Encryption Options	574
18.9. Secure TCP/IP Connections with SSL	575
18.10. Secure TCP/IP Connections with GSSAPI Encryption	579
18.11. Secure TCP/IP Connections with SSH Tunnels	579
18.12. Registering Event Log on Windows	580
19. Server Configuration	581
19.1. Setting Parameters	581
19.2. File Locations	584
19.3. Connections and Authentication	585
19.4. Resource Consumption	591
19.5. Write Ahead Log	602
19.6. Replication	612
19.7. Query Planning	619
19.8. Error Reporting and Logging	627
19.9. Run-time Statistics	641
19.10. Automatic Vacuuming	643
19.11. Background Freezer	645
19.12. Client Connection Defaults	645
19.13. Lock Management	656
19.14. Version and Platform Compatibility	657
19.15. Memory Purge	660
19.16. Data Compression	660
19.17. Error Handling	661
19.18. Preset Options	662
19.19. Customized Options	664
19.20. Developer Options	665
19.21. Distributed Cluster Configuration	670
19.22. Short Options	687
20. Client Authentication	689
20.1. The <code>pg_hba.conf</code> File	689
20.2. User Name Maps	697

20.3. Authentication Methods	698
20.4. Trust Authentication	699
20.5. Password Authentication	699
20.6. GSSAPI Authentication	700
20.7. SSPI Authentication	701
20.8. Ident Authentication	702
20.9. Peer Authentication	703
20.10. LDAP Authentication	703
20.11. RADIUS Authentication	706
20.12. Certificate Authentication	707
20.13. PAM Authentication	708
20.14. BSD Authentication	708
20.15. Authentication Problems	709
21. Database Roles	710
21.1. Database Roles	710
21.2. Role Attributes	711
21.3. Role Membership	712
21.4. Dropping Roles	714
21.5. Predefined Roles	714
21.6. Function Security	716
21.7. Roles in Distributed System	716
22. Managing Databases	718
22.1. Postgres Pro Shardman Overview	718
22.2. Connecting and Working with a Postgres Pro Shardman Cluster	718
22.3. Creating a Database	723
22.4. Template Databases	723
22.5. Database Configuration	724
22.6. Destroying a Database	725
22.7. Tablespaces	725
23. Localization	727
23.1. Locale Support	727
23.2. Collation Support	731
23.3. Character Set Support	740
24. Routine Database Maintenance Tasks	750
24.1. Routine Vacuuming	750
24.2. Background Freezer	758
24.3. Routine Reindexing	759
24.4. Log File Maintenance	759
25. Backup and Restore	761
25.1. Continuous Archiving and Point-in-Time Recovery (PITR)	761
25.2. Backup and Recovery in Distributed System	766
26. High Availability, Load Balancing, and Replication	770
26.1. Comparison of Different Solutions	770
26.2. Log-Shipping Standby Servers	773
26.3. Failover	782
26.4. Hot Standby	783
27. Monitoring Database Activity	791
27.1. Standard Unix Tools	791
27.2. The Cumulative Statistics System	792
27.3. Viewing Locks	834
27.4. Progress Reporting	835
27.5. Monitoring Disk Usage	842
28. Reliability and the Write-Ahead Log	845
28.1. Reliability	845
28.2. Data Checksums	846
28.3. Write-Ahead Logging (WAL)	847
28.4. Asynchronous Commit	847
28.5. WAL Configuration	849

28.6. WAL Restoration	852
28.7. WAL Internals	852
29. Logical Replication	854
29.1. Publication	854
29.2. Subscription	855
29.3. Logical Replication Failover	861
29.4. Row Filters	862
29.5. Column Lists	869
29.6. Conflicts	871
29.7. Restrictions	871
29.8. Architecture	872
29.9. Monitoring	873
29.10. Security	873
29.11. Configuration Settings	874
29.12. Quick Setup	875
30. Just-in-Time Compilation (JIT)	876
30.1. What Is JIT compilation?	876
30.2. When to JIT?	876
30.3. Configuration	877
30.4. Extensibility	878
31. Enhanced Security	879
31.1. Memory Purge	879
31.2. Integrity Checks	880
31.3. Separation of Duties between Privileged DBMS Users	882
31.4. Restricting DBMS Administrator's Data Access	885
32. Compressed File System (CFS)	890
32.1. Why database compression may be useful	890
32.2. How compression is integrated in Postgres Pro Shardman	890
32.3. Using Compression	891
33. Cluster Configuration for Distributed System	894
33.1. Cluster Initialization Settings Related to Access Management	894
33.2. Scaling the Cluster	895
33.3. Rebalancing the Data	897
33.4. Fault Tolerance and High Availability	902
33.5. Cluster Services	903
33.6. Configuring Secure Communications with etcd	903
34. Troubleshooting	907
IV. Client Interfaces	909
35. libpq — C Library	910
35.1. Database Connection Control Functions	910
35.2. Connection Status Functions	928
35.3. Command Execution Functions	934
35.4. Asynchronous Command Processing	948
35.5. Pipeline Mode	952
35.6. Retrieving Query Results in Chunks	956
35.7. Canceling Queries in Progress	957
35.8. The Fast-Path Interface	961
35.9. Asynchronous Notification	962
35.10. Functions Associated with the COPY Command	963
35.11. Control Functions	966
35.12. Miscellaneous Functions	968
35.13. Notice Processing	971
35.14. Event System	972
35.15. Environment Variables	978
35.16. The Password File	980
35.17. The Connection Service File	980
35.18. LDAP Lookup of Connection Parameters	981
35.19. SSL Support	981

35.20. Behavior in Threaded Programs	985
35.21. Building libpq Programs	985
35.22. Example Programs	987
36. Large Objects	997
36.1. Introduction	997
36.2. Implementation Features	997
36.3. Client Interfaces	997
36.4. Server-Side Functions	1001
36.5. Example Program	1002
37. ECPG — Embedded SQL in C	1008
37.1. The Concept	1008
37.2. Managing Database Connections	1008
37.3. Running SQL Commands	1011
37.4. Using Host Variables	1014
37.5. Dynamic SQL	1027
37.6. pgtypes Library	1028
37.7. Using Descriptor Areas	1040
37.8. Error Handling	1052
37.9. Preprocessor Directives	1058
37.10. Processing Embedded SQL Programs	1060
37.11. Library Functions	1061
37.12. Large Objects	1061
37.13. C++ Applications	1063
37.14. Embedded SQL Commands	1066
37.15. Informix Compatibility Mode	1088
37.16. Oracle Compatibility Mode	1101
37.17. Internals	1101
38. The Information Schema	1104
38.1. The Schema	1104
38.2. Data Types	1104
38.3. information_schema_catalog_name	1105
38.4. administrable_role_authorizations	1105
38.5. applicable_roles	1105
38.6. attributes	1106
38.7. character_sets	1108
38.8. check_constraint_routine_usage	1109
38.9. check_constraints	1109
38.10. collations	1109
38.11. collation_character_set_applicability	1110
38.12. column_column_usage	1110
38.13. column_domain_usage	1111
38.14. column_options	1111
38.15. column_privileges	1111
38.16. column_udt_usage	1112
38.17. columns	1113
38.18. constraint_column_usage	1116
38.19. constraint_table_usage	1116
38.20. data_type_privileges	1117
38.21. domain_constraints	1117
38.22. domain_udt_usage	1118
38.23. domains	1118
38.24. element_types	1120
38.25. enabled_roles	1122
38.26. foreign_data_wrapper_options	1122
38.27. foreign_data_wrappers	1123
38.28. foreign_server_options	1123
38.29. foreign_servers	1123

38.30. foreign_table_options	1124
38.31. foreign_tables	1124
38.32. key_column_usage	1125
38.33. parameters	1125
38.34. referential_constraints	1127
38.35. role_column_grants	1128
38.36. role_routine_grants	1128
38.37. role_table_grants	1129
38.38. role_udt_grants	1129
38.39. role_usage_grants	1130
38.40. routine_column_usage	1131
38.41. routine_privileges	1131
38.42. routine_routine_usage	1132
38.43. routine_sequence_usage	1132
38.44. routine_table_usage	1133
38.45. routines	1134
38.46. schemata	1138
38.47. sequences	1138
38.48. sql_features	1139
38.49. sql_implementation_info	1140
38.50. sql_parts	1140
38.51. sql_sizing	1140
38.52. table_constraints	1141
38.53. table_privileges	1141
38.54. tables	1142
38.55. transforms	1143
38.56. triggered_update_columns	1143
38.57. triggers	1144
38.58. udt_privileges	1145
38.59. usage_privileges	1146
38.60. user_defined_types	1146
38.61. user_mapping_options	1148
38.62. user_mappings	1148
38.63. view_column_usage	1149
38.64. view_routine_usage	1149
38.65. view_table_usage	1150
38.66. views	1150
V. Server Programming	1152
39. Extending SQL	1153
39.1. How Extensibility Works	1153
39.2. The Postgres Pro Type System	1153
39.3. User-Defined Functions	1156
39.4. User-Defined Procedures	1157
39.5. Query Language (SQL) Functions	1157
39.6. Function Overloading	1171
39.7. Function Volatility Categories	1172
39.8. Procedural Language Functions	1173
39.9. Internal Functions	1174
39.10. C-Language Functions	1174
39.11. Function Optimization Information	1194
39.12. User-Defined Aggregates	1195
39.13. User-Defined Types	1201
39.14. User-Defined Operators	1205
39.15. Operator Optimization Information	1206
39.16. Interfacing Extensions to Indexes	1209
39.17. Packaging Related Objects into an Extension	1221

39.18. Extension Building Infrastructure	1229
40. Triggers	1233
40.1. Overview of Trigger Behavior	1233
40.2. Visibility of Data Changes	1236
40.3. Writing Trigger Functions in C	1236
40.4. A Complete Trigger Example	1239
41. Event Triggers	1243
41.1. Overview of Event Trigger Behavior	1243
41.2. Event Trigger Firing Matrix	1244
41.3. Writing Event Trigger Functions in C	1247
41.4. A Complete Event Trigger Example	1248
41.5. A Table Rewrite Event Trigger Example	1249
41.6. A Database Login Event Trigger Example	1250
42. The Rule System	1252
42.1. The Query Tree	1252
42.2. Views and the Rule System	1253
42.3. Materialized Views	1260
42.4. Rules on INSERT, UPDATE, and DELETE	1262
42.5. Rules and Privileges	1271
42.6. Rules and Command Status	1273
42.7. Rules Versus Triggers	1274
43. Procedural Languages	1276
43.1. Installing Procedural Languages	1276
44. PL/pgSQL — SQL Procedural Language	1278
44.1. Overview	1278
44.2. Structure of PL/pgSQL	1279
44.3. Declarations	1280
44.4. Expressions	1286
44.5. Basic Statements	1287
44.6. Control Structures	1294
44.7. Cursors	1308
44.8. Transaction Management	1313
44.9. Errors and Messages	1314
44.10. Trigger Functions	1316
44.11. Packages	1324
44.12. PL/pgSQL under the Hood	1325
44.13. Tips for Developing in PL/pgSQL	1328
44.14. Porting from Oracle PL/SQL	1331
45. PL/Tcl — Tcl Procedural Language	1340
45.1. Overview	1340
45.2. PL/Tcl Functions and Arguments	1340
45.3. Data Values in PL/Tcl	1342
45.4. Global Data in PL/Tcl	1342
45.5. Database Access from PL/Tcl	1342
45.6. Trigger Functions in PL/Tcl	1344
45.7. Event Trigger Functions in PL/Tcl	1346
45.8. Error Handling in PL/Tcl	1347
45.9. Explicit Subtransactions in PL/Tcl	1347
45.10. Transaction Management	1348
45.11. PL/Tcl Configuration	1349
45.12. Tcl Procedure Names	1349
46. PL/Perl — Perl Procedural Language	1350
46.1. PL/Perl Functions and Arguments	1350
46.2. Data Values in PL/Perl	1354
46.3. Built-in Functions	1354
46.4. Global Values in PL/Perl	1359
46.5. Trusted and Untrusted PL/Perl	1359
46.6. PL/Perl Triggers	1361

46.7. PL/Perl Event Triggers	1362
46.8. PL/Perl Under the Hood	1362
47. PL/Python — Python Procedural Language	1365
47.1. PL/Python Functions	1365
47.2. Data Values	1366
47.3. Sharing Data	1371
47.4. Anonymous Code Blocks	1371
47.5. Trigger Functions	1371
47.6. Database Access	1372
47.7. Explicit Subtransactions	1375
47.8. Transaction Management	1376
47.9. Utility Functions	1376
47.10. Python 2 vs. Python 3	1377
47.11. Environment Variables	1377
48. Server Programming Interface	1379
48.1. Interface Functions	1379
48.2. Interface Support Functions	1418
48.3. Memory Management	1427
48.4. Transaction Management	1437
48.5. Visibility of Data Changes	1440
48.6. Examples	1440
49. Background Worker Processes	1444
50. Logical Decoding	1447
50.1. Logical Decoding Examples	1447
50.2. Logical Decoding Concepts	1450
50.3. Streaming Replication Protocol Interface	1452
50.4. Logical Decoding SQL Interface	1452
50.5. System Catalogs Related to Logical Decoding	1453
50.6. Logical Decoding Output Plugins	1453
50.7. Logical Decoding Output Writers	1459
50.8. Synchronous Replication Support for Logical Decoding	1459
50.9. Streaming of Large Transactions for Logical Decoding	1460
50.10. Two-phase Commit Support for Logical Decoding	1461
51. Replication Progress Tracking	1462
52. Archive Modules	1463
52.1. Initialization Functions	1463
52.2. Archive Module Callbacks	1463
VI. Reference	1465
I. SQL Commands	1466
ABORT	1467
ALTER AGGREGATE	1468
ALTER COLLATION	1470
ALTER CONVERSION	1472
ALTER DATABASE	1473
ALTER DEFAULT PRIVILEGES	1475
ALTER DOMAIN	1478
ALTER EVENT TRIGGER	1481
ALTER EXTENSION	1482
ALTER FUNCTION	1485
ALTER GROUP	1488
ALTER INDEX	1490
ALTER LANGUAGE	1493
ALTER LARGE OBJECT	1494
ALTER MATERIALIZED VIEW	1495
ALTER OPERATOR	1497
ALTER OPERATOR CLASS	1499
ALTER OPERATOR FAMILY	1500
ALTER POLICY	1504

ALTER PROCEDURE	1505
ALTER PROFILE	1508
ALTER PUBLICATION	1510
ALTER ROLE	1512
ALTER ROUTINE	1516
ALTER RULE	1517
ALTER SCHEMA	1518
ALTER SEQUENCE	1519
ALTER SERVER	1523
ALTER STATISTICS	1524
ALTER SUBSCRIPTION	1525
ALTER TABLE	1528
ALTER TABLESPACE	1548
ALTER TEXT SEARCH CONFIGURATION	1549
ALTER TEXT SEARCH DICTIONARY	1551
ALTER TEXT SEARCH PARSER	1553
ALTER TEXT SEARCH TEMPLATE	1554
ALTER TRIGGER	1555
ALTER TYPE	1557
ALTER USER	1561
ALTER VIEW	1562
ANALYZE	1564
BEGIN	1567
CALL	1569
CHECKPOINT	1570
CLOSE	1571
CLUSTER	1572
COMMENT	1575
COMMIT	1579
COMMIT PREPARED	1580
COPY	1581
CREATE ACCESS METHOD	1591
CREATE AGGREGATE	1592
CREATE CAST	1599
CREATE COLLATION	1603
CREATE CONVERSION	1606
CREATE DATABASE	1608
CREATE DOMAIN	1613
CREATE EVENT TRIGGER	1616
CREATE EXTENSION	1618
CREATE FUNCTION	1620
CREATE GROUP	1628
CREATE INDEX	1629
CREATE LANGUAGE	1637
CREATE MATERIALIZED VIEW	1639
CREATE OPERATOR	1641
CREATE OPERATOR CLASS	1644
CREATE OPERATOR FAMILY	1647
CREATE PACKAGE	1648
CREATE POLICY	1652
CREATE PROCEDURE	1657
CREATE PROFILE	1661
CREATE PUBLICATION	1665
CREATE ROLE	1669
CREATE RULE	1674
CREATE SCHEMA	1677
CREATE SEQUENCE	1679
CREATE SERVER	1683

CREATE STATISTICS	1685
CREATE SUBSCRIPTION	1689
CREATE TABLE	1694
CREATE TABLE AS	1719
CREATE TABLESPACE	1722
CREATE TEXT SEARCH CONFIGURATION	1724
CREATE TEXT SEARCH DICTIONARY	1725
CREATE TEXT SEARCH PARSER	1727
CREATE TEXT SEARCH TEMPLATE	1729
CREATE TRANSFORM	1730
CREATE TRIGGER	1732
CREATE TYPE	1739
CREATE USER	1748
CREATE VIEW	1749
DEALLOCATE	1754
DECLARE	1755
DELETE	1758
DISCARD	1761
DO	1762
DROP ACCESS METHOD	1763
DROP AGGREGATE	1764
DROP CAST	1766
DROP COLLATION	1767
DROP CONVERSION	1768
DROP DATABASE	1769
DROP DOMAIN	1770
DROP EVENT TRIGGER	1771
DROP EXTENSION	1772
DROP FUNCTION	1773
DROP GROUP	1775
DROP INDEX	1776
DROP LANGUAGE	1777
DROP MATERIALIZED VIEW	1778
DROP OPERATOR	1779
DROP OPERATOR CLASS	1780
DROP OPERATOR FAMILY	1781
DROP OWNED	1782
DROP PACKAGE	1783
DROP POLICY	1784
DROP PROCEDURE	1785
DROP PROFILE	1787
DROP PUBLICATION	1788
DROP ROLE	1789
DROP ROUTINE	1790
DROP RULE	1791
DROP SCHEMA	1792
DROP SEQUENCE	1793
DROP SERVER	1794
DROP STATISTICS	1795
DROP SUBSCRIPTION	1796
DROP TABLE	1797
DROP TABLESPACE	1798
DROP TEXT SEARCH CONFIGURATION	1799
DROP TEXT SEARCH DICTIONARY	1800
DROP TEXT SEARCH PARSER	1801
DROP TEXT SEARCH TEMPLATE	1802
DROP TRANSFORM	1803
DROP TRIGGER	1804

DROP TYPE	1805
DROP USER	1806
DROP VIEW	1807
END	1808
EXECUTE	1809
EXPLAIN	1810
FETCH	1816
GRANT	1820
INSERT	1826
LISTEN	1833
LOAD	1835
LOCK	1836
MERGE	1839
MOVE	1846
NOTIFY	1848
PREPARE	1850
PREPARE TRANSACTION	1853
REASSIGN OWNED	1855
REFRESH MATERIALIZED VIEW	1856
REINDEX	1858
RELEASE SAVEPOINT	1863
RESET	1865
REVOKE	1866
ROLLBACK	1870
ROLLBACK PREPARED	1871
ROLLBACK TO SAVEPOINT	1872
SAVEPOINT	1874
SECURITY LABEL	1876
SELECT	1878
SELECT INTO	1898
SET	1900
SET CONSTRAINTS	1903
SET ROLE	1904
SET SESSION AUTHORIZATION	1906
SET TRANSACTION	1908
SHOW	1911
START TRANSACTION	1913
TRUNCATE	1914
UNLISTEN	1916
UPDATE	1917
VACUUM	1922
VALUES	1927
II. Postgres Pro Client Applications	1929
clusterdb	1930
createdb	1933
createuser	1937
dropdb	1942
dropuser	1945
ecpg	1947
pg_amcheck	1949
pgbench	1954
pg_combinebackup	1977
pg_config	1980
pg_dump	1983
pg_dumpall	1997
pg_isready	2004
pg_receivewal	2006
pg_recvlogical	2010

pg_restore	2014
pg-wrapper	2023
pg_verifybackup	2025
psql	2028
reindexdb	2070
vacuumdb	2073
III. Postgres Pro Server Applications	2078
initdb	2079
pg_archivecleanup	2085
pg_checksums	2087
pg_createsubscriber	2089
pg_ctl	2093
pgpro_tune	2098
pg_resetwal	2102
pg_rewind	2106
pg-setup	2110
pg_test_fsync	2112
pg_test_timing	2113
pg_waldump	2116
pg_walsummary	2120
postgres	2121
shardmanctl	2128
shardmand	2194
VII. Internals	2197
53. Overview of Postgres Pro Internals	2198
53.1. The Path of a Query	2198
53.2. Query Processing	2198
53.3. How Connections Are Established	2201
53.4. The Parser Stage	2202
53.5. The Postgres Pro Rule System	2203
53.6. Planner/Optimizer	2203
53.7. Executor	2204
54. System Catalogs	2206
54.1. Overview	2206
54.2. pg_aggregate	2208
54.3. pg_am	2209
54.4. pg_amop	2210
54.5. pg_amproc	2211
54.6. pg_attrdef	2211
54.7. pg_attribute	2211
54.8. pg_authid	2213
54.9. pg_auth_members	2215
54.10. pg_cast	2215
54.11. pg_class	2216
54.12. pg_collation	2218
54.13. pg_constraint	2219
54.14. pg_conversion	2221
54.15. pg_database	2222
54.16. pg_db_role_setting	2223
54.17. pg_default_acl	2223
54.18. pg_depend	2224
54.19. pg_description	2226
54.20. pg_enum	2226
54.21. pg_event_trigger	2227
54.22. pg_extension	2227
54.23. pg_foreign_data_wrapper	2228
54.24. pg_foreign_server	2228

54.25. pg_foreign_table	2229
54.26. pg_index	2229
54.27. pg_inherits	2231
54.28. pg_init_privs	2231
54.29. pg_language	2232
54.30. pg_largeobject	2232
54.31. pg_largeobject_metadata	2233
54.32. pg_namespace	2233
54.33. pg_opclass	2234
54.34. pg_operator	2234
54.35. pg_opfamily	2235
54.36. pg_parameter_acl	2236
54.37. pg_partitioned_table	2236
54.38. pg_policy	2237
54.39. pg_proc	2237
54.40. pg_profile	2240
54.41. pg_publication	2240
54.42. pg_publication_namespace	2241
54.43. pg_publication_rel	2241
54.44. pg_range	2242
54.45. pg_replication_origin	2242
54.46. pg_rewrite	2243
54.47. pg_role_password	2243
54.48. pg_seclabel	2243
54.49. pg_sequence	2244
54.50. pg_shdepend	2244
54.51. pg_shdescription	2246
54.52. pg_shseclabel	2246
54.53. pg_statistic	2246
54.54. pg_statistic_ext	2248
54.55. pg_statistic_ext_data	2248
54.56. pg_subscription	2249
54.57. pg_subscription_rel	2250
54.58. pg_tablespace	2251
54.59. pg_transform	2251
54.60. pg_trigger	2252
54.61. pg_ts_config	2253
54.62. pg_ts_config_map	2254
54.63. pg_ts_dict	2254
54.64. pg_ts_parser	2254
54.65. pg_ts_template	2255
54.66. pg_type	2255
54.67. pg_user_mapping	2259
55. System Views	2260
55.1. Overview	2260
55.2. pg_available_extensions	2261
55.3. pg_available_extension_versions	2261
55.4. pg_backend_memory_contexts	2262
55.5. pg_config	2263
55.6. pg_cursors	2263
55.7. pg_file_settings	2264
55.8. pg_group	2264
55.9. pg_hba_file_rules	2265
55.10. pg_ident_file_mappings	2265
55.11. pg_indexes	2266
55.12. pg_locks	2266

55.13.	pg_matviews	2269
55.14.	pg_policies	2269
55.15.	pg_prepared_statements	2270
55.16.	pg_autoprepared_statements	2271
55.17.	pg_prepared_xacts	2271
55.18.	pg_publication_tables	2272
55.19.	pg_replication_origin_status	2272
55.20.	pg_replication_slots	2273
55.21.	pg_roles	2274
55.22.	pg_rules	2276
55.23.	pg_seclabels	2276
55.24.	pg_sequences	2276
55.25.	pg_settings	2277
55.26.	pg_shadow	2279
55.27.	pg_shmem_allocations	2280
55.28.	pg_stats	2280
55.29.	pg_stats_ext	2282
55.30.	pg_stats_ext_exprs	2283
55.31.	pg_stats_vacuum_database	2285
55.32.	pg_stats_vacuum_indexes	2286
55.33.	pg_stats_vacuum_tables	2288
55.34.	pg_tables	2289
55.35.	pg_timezone_abbrevs	2290
55.36.	pg_timezone_names	2290
55.37.	pg_user	2290
55.38.	pg_user_mappings	2291
55.39.	pg_views	2292
55.40.	pg_wait_events	2292
55.41.	Global Views for Distributed System	2292
56.	Frontend/Backend Protocol	2301
56.1.	Overview	2301
56.2.	Message Flow	2302
56.3.	SASL Authentication	2315
56.4.	Streaming Replication Protocol	2317
56.5.	Logical Streaming Replication Protocol	2326
56.6.	Message Data Types	2327
56.7.	Message Formats	2328
56.8.	Error and Notice Message Fields	2343
56.9.	Logical Replication Message Formats	2345
56.10.	Summary of Changes since Protocol 2.0	2353
57.	Writing a Procedural Language Handler	2354
58.	Writing a Foreign Data Wrapper	2356
58.1.	Foreign Data Wrapper Functions	2356
58.2.	Foreign Data Wrapper Callback Routines	2356
58.3.	Foreign Data Wrapper Helper Functions	2370
58.4.	Foreign Data Wrapper Query Planning	2371
58.5.	Row Locking in Foreign Data Wrappers	2373
59.	Writing a Table Sampling Method	2375
59.1.	Sampling Method Support Functions	2375
60.	Writing a Custom Scan Provider	2378
60.1.	Creating Custom Scan Paths	2378
60.2.	Creating Custom Scan Plans	2379
60.3.	Executing Custom Scans	2380
61.	Genetic Query Optimizer	2383
61.1.	Query Handling as a Complex Optimization Problem	2383
61.2.	Genetic Algorithms	2383
61.3.	Genetic Query Optimization (GEQO) in Postgres Pro	2384

61.4. Further Reading	2386
62. Table Access Method Interface Definition	2387
63. Index Access Method Interface Definition	2388
63.1. Basic API Structure for Indexes	2388
63.2. Index Access Method Functions	2391
63.3. Index Scanning	2396
63.4. Index Locking Considerations	2398
63.5. Index Uniqueness Checks	2399
63.6. Index Cost Estimation Functions	2400
64. Write Ahead Logging for Extensions	2403
64.1. Generic WAL Records	2403
64.2. Custom WAL Resource Managers	2404
65. Built-in Index Access Methods	2406
65.1. B-Tree Indexes	2406
65.2. GiST Indexes	2412
65.3. SP-GiST Indexes	2427
65.4. GIN Indexes	2438
65.5. BRIN Indexes	2444
65.6. Hash Indexes	2456
66. Database Physical Storage	2459
66.1. Database File Layout	2459
66.2. TOAST	2461
66.3. Free Space Map	2466
66.4. Visibility Map	2467
66.5. The Initialization Fork	2467
66.6. Database Page Layout	2467
66.7. Heap-Only Tuples (HOT)	2470
67. Transaction Processing	2471
67.1. Transactions and Identifiers	2471
67.2. Transactions and Locking	2471
67.3. Subtransactions	2471
67.4. Two-Phase Transactions	2472
67.5. Distributed Transactions	2472
68. How the Planner Uses Statistics	2477
68.1. Row Estimation Examples	2477
68.2. Multivariate Statistics Examples	2481
68.3. Planner Statistics and Security	2484
69. Backup Manifest Format	2485
69.1. Backup Manifest Top-level Object	2485
69.2. Backup Manifest File Object	2485
69.3. Backup Manifest WAL Range Object	2486
70. Distributed System Internals	2487
70.1. Table Types	2487
70.2. Silk	2488
70.3. Distributed Deadlock Detection	2500
70.4. Global Sequences	2501
70.5. Syncpoints and Consistent Backup	2502
70.6. Collecting Distributed Statement Statistics Using the pgpro_stats Extension	2502
70.7. Advisory Locks	2503
VIII. Appendixes	2504
A. Postgres Pro Error Codes	2505
B. Date/Time Support	2514
B.1. Date/Time Input Interpretation	2514
B.2. Handling of Invalid or Ambiguous Timestamps	2515
B.3. Date/Time Key Words	2515
B.4. Date/Time Configuration Files	2516
B.5. POSIX Time Zone Specifications	2518
B.6. History of Units	2519

B.7. Julian Dates	2520
C. SQL Key Words	2522
D. SQL Conformance	2548
D.1. Supported Features	2549
D.2. Unsupported Features	2561
D.3. XML Limits and Conformance to SQL/XML	2569
E. Release Notes	2572
E.1. Postgres Pro Shardman 17.5.1	2572
E.2. Release 17.5	2574
E.3. Release 17.4	2579
E.4. Release 17.3	2579
E.5. Release 17.2	2585
E.6. Release 17.1	2586
E.7. Release 17	2591
E.8. Prior Releases	2610
F. Additional Supplied Modules and Extensions Shipped in <code>postgrespro-sdm-17-contrib</code>	2611
F.1. <code>amcheck</code> — tools to verify table and index consistency	2612
F.2. <code>aqo</code> — cost-based query optimization	2617
F.3. <code>auto_explain</code> — log execution plans of slow queries	2639
F.4. <code>dblink</code> — connect to other Postgres Pro databases	2642
F.5. <code>pgcrypto</code> — cryptographic functions	2670
F.6. <code>pg_proaudit</code> — enables detailed logging of various security events	2679
F.7. <code>pgpro_stats</code> — a means for tracking planning and execution statistics of all SQL state- ments executed by a server	2692
F.8. <code>pg_query_state</code> — a facility to know the current state of query execution on working backend	2734
F.9. <code>pg_stat_statements</code> — track statistics of SQL planning and execution	2740
F.10. <code>pg_wait_sampling</code> — collecting sampling-based statistics on wait events	2748
F.11. <code>postgres_fdw</code> — access data stored in external Postgres Pro servers	2752
F.12. <code>ptrack</code> — a block-level incremental backup engine for Postgres Pro	2762
F.13. <code>uuid-osp</code> — a UUID generator	2764
G. Postgres Pro Modules and Extensions Shipped as Individual Packages	2766
G.1. <code>pgpro_pwr</code> — workload reports	2766
H. Additional Supplied Programs	2831
H.1. Additional PostgreSQL/Postgres Pro Client Applications	2831
H.2. Additional Postgres Pro Server Applications	2900
H.3. Third-Party Server Applications	2904
I. External Projects	2938
I.1. Client Interfaces	2938
I.2. Administration Tools	2938
I.3. Procedural Languages	2938
I.4. Extensions	2938
J. Postgres Pro Limits	2939
K. Demo Database “Airlines”	2940
K.1. Installation	2940
K.2. Schema Diagram	2941
K.3. Schema Description	2941
K.4. Schema Objects	2942
K.5. Usage	2948
L. Acronyms	2956
M. Glossary	2961
N. Color Support	2974
N.1. When Color is Used	2974
N.2. Configuring the Colors	2974
O. Obsolete or Renamed Features	2975
O.1. <code>recovery.conf</code> file merged into <code>postgresql.conf</code>	2975
O.2. Default Roles Renamed to Predefined Roles	2975
O.3. <code>pg_xlogdump</code> renamed to <code>pg_waldump</code>	2975

O.4. pg_resetxlog renamed to pg_resetwal	2975
O.5. pg_receivexlog renamed to pg_receivewal	2975
Bibliography	2976
Index	2978

Preface

This book is the official documentation of Postgres Pro Shardman. It has been written by the Postgres Pro developers, PostgreSQL community, and other volunteers in parallel to the development of the PostgreSQL and Postgres Pro software. It describes all the functionality that the current version of Postgres Pro officially supports.

To make the large amount of information about Postgres Pro manageable, this book has been organized in several parts. Each part is targeted at a different class of users, or at users in different stages of their Postgres Pro experience:

- [Part I](#) is an informal introduction for new users.
- [Part II](#) documents the SQL query language environment, including data types and functions, as well as user-level performance tuning. Every Postgres Pro user should read this.
- [Part III](#) describes the installation and administration of the server. Everyone who runs a Postgres Pro server, be it for private use or for others, should read this part.
- [Part IV](#) describes the programming interfaces for Postgres Pro client programs.
- [Part V](#) contains information for advanced users about the extensibility capabilities of the server. Topics include user-defined data types and functions.
- [Part VI](#) contains reference information about SQL commands, client and server programs. This part supports the other parts with structured information sorted by command or program.
- [Part VII](#) contains assorted information that might be of use to Postgres Pro developers.

1. What Is Postgres Pro Shardman?

Postgres Pro Shardman is based on the Postgres Pro Enterprise and is an object-relational database management system (ORDBMS), developed by Postgres Professional in the Postgres Pro fork of *PostgreSQL*, which is in turn based on *POSTGRES, Version 4.2*, developed at the University of California at Berkeley Computer Science Department. POSTGRES pioneered many concepts that only became available in some commercial database systems much later.

The database size in modern enterprises and in highload web applications is constantly growing. The only working approach to accommodate this growth is horizontal scaling. The Postgres Pro Shardman distributed DBMS is intended to enable horizontal scaling of online transaction processing (OLTP) databases while preserving the strong [ACID](#) semantics.

Postgres Pro Shardman runs on all major [Linux](#) operating systems.

Both PostgreSQL and Postgres Pro Shardman support a large part of the SQL standard and offer many modern features:

- [complex queries](#)
- [foreign keys](#)
- [triggers](#)
- [updatable views](#)
- [transactional integrity](#)
- [multiversion concurrency control](#)

Besides, PostgreSQL and Postgres Pro Shardman can be extended by the user in many ways, for example by adding new

- [data types](#)
- [functions](#)
- [operators](#)
- [aggregate functions](#)
- [index methods](#)
- [procedural languages](#)

Postgres Pro Shardman also provides the following advantages, compatibility features to your applications:

- Strong ACID guarantees.
- Compatibility with *Postgres Pro Enterprise*.
- Trust level 4 and security class 4 certificates.
- Several clusters support.
- Transparent [horizontal scaling](#) without a need in adopting NoSQL DBMS.
- [Built-in support of replication with no single point of failure](#), with any node being able to become coordinator that requires no system shut down and prevents any data loss.
- Capacity of up to 100 cluster nodes.
- [High availability](#) with primary and stand-by modes, along with the synchronous solution for the failover scenarios, and asynchronous solution that has a minimal performance impact.
- [Support](#) of planning and execution statistics of all SQL statements.
- [Utility](#) to discover most resource-intensive activities in your database.
- [Tools](#) to support `REPEATABLE READ` isolation level in a distributed system.
- [Work with cluster](#) as with a fully functional DBMS.
- [Hot standby](#) and [backup and recovery](#) tools that support full and incremental backup with logs.
- [Point-in-time recovery \(PITR\)](#).
- [Streaming replication](#).
- High-availability cluster creation with multiple primary nodes with special [utilities](#).
- [ANSI](#) standard.
- [Stored procedures](#).
- [Big data](#) storing and processing.
- [Full text search](#).
- [Covering indexes](#).
- [B-tree](#), [hash](#), [GiST](#), [GIN](#), [SP-GiST](#), [BRIN](#) indexes.
- [Perl](#) and [Python](#) procedural languages.
- Interfaces for [C++](#), [Ruby](#), [C](#), [ODBC](#), [Perl](#), [Python](#), [Tcl](#), and [Java](#).
- [EUC](#), [UTF-8](#), and [Mule](#) character set.
- [Adaptive query optimization](#).
- [Compressed file system](#).
- Access data stored in external PostgreSQL servers with [postgres_fdw](#), e.g. [Microsoft Active Directory](#), [Mysql server](#), [Oracle](#), and *Postgres Pro Enterprise*.
- Long queries monitoring with the [pg_query_state module](#).
- [In-built monitoring agent](#).
- [No limits](#) for the number of records or indexes, with the maximum table size of 32 TB, maximum attribute size of 1 GB, and maximum number of attributes of 1600.
- [Detailed access management](#) with different access levels and roles.
- [Secure password storing](#).
- Detailed [memory purge](#) configuration.

2. Difference between Postgres Pro Shardman and PostgreSQL

Postgres Pro provides the most actual PostgreSQL version with some additional patches applied and extensions added. It includes new features developed by Postgres Professional, as well as third-party patches already accepted by the PostgreSQL community for the upcoming PostgreSQL versions. Postgres Pro Shardman users thus have early access to important features and fixes.

Note

Postgres Pro Shardman is provided under the following license: <https://postgrespro.com/products/postgrespro/eula>. Make sure to review the license terms before downloading Postgres Pro Shardman.

Postgres Pro Shardman provides the following enhancements over PostgreSQL:

- 64-bit transaction IDs that are not subject to wraparound. (See [Section 24.1.5](#).)
- Page-level compression. (See [Chapter 32](#).)
- Support for autonomous transactions. (See [Chapter 16](#).)
- Fair lightweight lock scheduling after the specified number of shared locks is acquired. (See [lwlock_shared_limit](#) parameter description.)
- Improved multi-host connection handling and failover by libpq. (See [hostorder](#) and [failover_timeout](#) parameter descriptions.)
- Enabling libpq to forget the entered password, which allows to prevent reconnections when required by a security policy. (See [reusepass](#) parameter description.)
- Support for timestamp output in [pg_waldump](#).
- K-nearest neighbors (k-NN) algorithm for B-tree indexes. (See [Section 11.13](#).)
- Removing a practically reachable limit on the number of entries in ACL (access control list, i.e., privileges list) associated with tables and indexes.
- The [pgpro_stat_wal_activity](#) view that shows the size of WAL files generated by each process.
- Verification of unique constraints in B-tree indexes in the [amcheck](#) module.
- Support for passing named and positional arguments to scripts invoked by the `\i` command in [psql](#).
- [pg_probackup](#) enterprise edition, which provides Simple Storage Service (S3) support for storing data in private clouds, [CFS](#) (Compressed File System) support for incremental backups, and support for lz4 and zstd compression algorithms.
- Lock deduplication that allows to effectively store in memory and track all exclusive locks held by a standby server's startup process during WAL replay.
- Optimized mechanism for working with table metadata, which allows obtaining information about attributes using the system cache instead of direct reading from the system catalog.

The following enhancements are inherited from Postgres Pro Standard:

- Better planning speed and accuracy for various query types.
- Reduced memory consumption in complex queries that involve multiple tables.
- Displaying planning time in the output of the [auto_explain](#) module.
- NUL byte replacement with the specified ASCII code while loading data using the `COPY FROM` command. (See [nul_byte_replacement_on_import](#) parameter description.)
- `'\u0000'` character replacement with the specified unicode character when calling a function processing JSONB (See [unicode_nul_character_replacement_in_jsonb](#) parameter description.)
- [PTRACK](#) implementation, which enables [pg_probackup](#) to track page changes on the fly when creating incremental backups.
- Changing the `restore_command` parameter without restarting the server.
- Advanced authentication policies that provide effective password management and access control. (See [CREATE PROFILE](#) and [ALTER ROLE](#)).
- Built-in [data security mechanisms](#) that enable sanitizing an object by filling it with zeroes before deletion. Zeroing can be done before purging files in external memory and removing outdated row

versions (page vacuum), freeing RAM, and deleting or overwriting WAL files. (Certified edition only.)

- Statistics about vacuuming tables, indexes and databases in [system views](#).
- [Predefined roles](#), which allow creating tablespaces and managing profiles without superuser rights.
- Getting information on crashes of a backend, which is enabled by the [crash_info](#) configuration parameter and controlled by more of them.
- Optimized memory consumption during selectivity estimation for each array element.

Postgres Pro Shardman also includes the following additional modules and applications:

- [aqo](#) extension for adaptive query optimization.
- [pgbouncer](#) connection pooler.
- [pg_integrity_check](#) module that calculates and validates checksums for controlled files. (Certified edition only.)
- [pg_proaudit](#) extension that enables detailed logging of various security events.
- [pg_probackup](#), a backup and recovery manager.
- [pgpro_controldata](#), an application to display control information of a PostgreSQL/Postgres Pro database cluster and compatibility information for a cluster and/or server.
- [pgpro_pwr](#) extension that enables you to generate workload reports, which help to discover most resource-intensive activities in your database.
- [pgpro_stats](#) extension that tracks execution statistics of SQL statements, calculates wait event statistics and provides other useful metrics that are not collected elsewhere in PostgreSQL. It also provides tracing of application sessions and can create views that emulate other statistic collecting extensions.
- [pg_query_state](#) module that enables you to get the current state of query execution for a backend.
- [pgvector](#) extension that provides vector similarity search for Postgres Pro.

To provide the advanced functionalities and features, Postgres Pro imposes more stringent requirements on operating systems supported.

Important

Postgres Pro Shardman runs on all major Linux operating systems. Any reference to Windows or another operating system different from Linux in this documentation is inapplicable for Postgres Pro.

Postgres Pro Shardman releases follow PostgreSQL releases, though sometimes occur more frequently. The Postgres Pro Shardman versioning scheme is based on the PostgreSQL one and has an additional decimal place.

3. A Brief History of PostgreSQL

The object-relational database management system now known as PostgreSQL is derived from the POSTGRES package written at the University of California at Berkeley. With decades of development behind it, PostgreSQL is now the most advanced open-source database available anywhere.

3.1. The Berkeley POSTGRES Project

The POSTGRES project, led by Professor Michael Stonebraker, was sponsored by the Defense Advanced Research Projects Agency (DARPA), the Army Research Office (ARO), the National Science Foundation (NSF), and ESL, Inc. The implementation of POSTGRES began in 1986. The initial concepts for the system were presented in [ston86](#), and the definition of the initial data model appeared in [rowe87](#). The design of the rule system at that time was described in [ston87a](#). The rationale and architecture of the storage manager were detailed in [ston87b](#).

POSTGRES has undergone several major releases since then. The first “demoware” system became operational in 1987 and was shown at the 1988 ACM-SIGMOD Conference. Version 1, described in [ston90a](#), was released to a few external users in June 1989. In response to a critique of the first rule

system ([ston89](#)), the rule system was redesigned ([ston90b](#)), and Version 2 was released in June 1990 with the new rule system. Version 3 appeared in 1991 and added support for multiple storage managers, an improved query executor, and a rewritten rule system. For the most part, subsequent releases until Postgres95 (see below) focused on portability and reliability.

POSTGRES has been used to implement many different research and production applications. These include: a financial data analysis system, a jet engine performance monitoring package, an asteroid tracking database, a medical information database, and several geographic information systems. POSTGRES has also been used as an educational tool at several universities. Finally, Illustra Information Technologies (later merged into *Informix*, which is now owned by *IBM*) picked up the code and commercialized it. In late 1992, POSTGRES became the primary data manager for the Sequoia 2000 scientific computing project described in [ston92](#).

The size of the external user community nearly doubled during 1993. It became increasingly obvious that maintenance of the prototype code and support was taking up large amounts of time that should have been devoted to database research. In an effort to reduce this support burden, the Berkeley POSTGRES project officially ended with Version 4.2.

3.2. Postgres95

In 1994, Andrew Yu and Jolly Chen added an SQL language interpreter to POSTGRES. Under a new name, Postgres95 was subsequently released to the web to find its own way in the world as an open-source descendant of the original POSTGRES Berkeley code.

Postgres95 code was completely ANSI C and trimmed in size by 25%. Many internal changes improved performance and maintainability. Postgres95 release 1.0.x ran about 30-50% faster on the Wisconsin Benchmark compared to POSTGRES, Version 4.2. Apart from bug fixes, the following were the major enhancements:

- The query language PostQUEL was replaced with SQL (implemented in the server). (Interface library [libpq](#) was named after PostQUEL.) Subqueries were not supported until PostgreSQL (see below), but they could be imitated in Postgres95 with user-defined SQL functions. Aggregate functions were re-implemented. Support for the `GROUP BY` query clause was also added.
- A new program (`psql`) was provided for interactive SQL queries, which used GNU Readline. This largely superseded the old monitor program.
- A new front-end library, `libpgtcl`, supported Tcl-based clients. A sample shell, `pgtclsh`, provided new Tcl commands to interface Tcl programs with the Postgres95 server.
- The large-object interface was overhauled. The inversion large objects were the only mechanism for storing large objects. (The inversion file system was removed.)
- The instance-level rule system was removed. Rules were still available as rewrite rules.
- A short tutorial introducing regular SQL features as well as those of Postgres95 was distributed with the source code
- GNU make (instead of BSD make) was used for the build. Also, Postgres95 could be compiled with an unpatched GCC (data alignment of doubles was fixed).

3.3. PostgreSQL

By 1996, it became clear that the name “Postgres95” would not stand the test of time. We chose a new name, PostgreSQL, to reflect the relationship between the original POSTGRES and the more recent versions with SQL capability. At the same time, we set the version numbering to start at 6.0, putting the numbers back into the sequence originally begun by the Berkeley POSTGRES project.

Postgres is still considered an official project name, both because of tradition and because people find it easier to pronounce Postgres than PostgreSQL.

The emphasis during development of Postgres95 was on identifying and understanding existing problems in the server code. With PostgreSQL, the emphasis has shifted to augmenting features and capabilities, although work continues in all areas.

4. Conventions

The following conventions are used in the synopsis of a command: brackets (`[` and `]`) indicate optional parts. Braces (`{` and `}`) and vertical lines (`|`) indicate that you must choose one alternative. Dots (`...`) mean that the preceding element can be repeated. All other symbols, including parentheses, should be taken literally.

Where it enhances the clarity, SQL commands are preceded by the prompt `=>`, and shell commands are preceded by the prompt `$`. Normally, prompts are not shown, though.

An *administrator* is generally a person who is in charge of installing and running the server. A *user* could be anyone who is using, or wants to use, any part of the Postgres Pro system. These terms should not be interpreted too narrowly; this book does not have fixed presumptions about system administration procedures.

5. Bug Reporting Guidelines

When you find a bug in Postgres Pro we want to hear about it. Your bug reports play an important part in making Postgres Pro more reliable because even the utmost care cannot guarantee that every part of Postgres Pro will work on every platform under every circumstance.

The following suggestions are intended to assist you in forming bug reports that can be handled in an effective fashion. No one is required to follow them but doing so tends to be to everyone's advantage.

5.1. Identifying Bugs

Before you report a bug, please read and re-read the documentation to verify that you can really do whatever it is you are trying. If it is not clear from the documentation whether you can do something or not, please report that too; it is a bug in the documentation. If it turns out that a program does something different from what the documentation says, that is a bug. That might include, but is not limited to, the following circumstances:

- A program terminates with a fatal signal or an operating system error message that would point to a problem in the program. (A counterexample might be a “disk full” message, since you have to fix that yourself.)
- A program produces the wrong output for any given input.
- A program refuses to accept valid input (as defined in the documentation).
- A program accepts invalid input without a notice or error message. But keep in mind that your idea of invalid input might be our idea of an extension or compatibility with traditional practice.
- Postgres Pro fails to install according to the instructions on supported platforms.

Here “program” refers to any executable, not only the backend process.

Being slow or resource-hogging is not necessarily a bug. Failing to comply to the SQL standard is not necessarily a bug either, unless compliance for the specific feature is explicitly claimed.

5.2. What to Report

When reporting a bug, make sure to state all the facts. Each bug report should contain the following items:

- The exact sequence of steps *from program start-up* necessary to reproduce the problem. This should be self-contained; it is not enough to send in a bare `SELECT` statement without the preceding `CREATE TABLE` and `INSERT` statements, if the output should depend on the data in the tables.

The best format for a test case for SQL-related problems is a file that can be run through the `psql` frontend that shows the problem. (Be sure to not have anything in your `~/.psqlrc` start-up file.) An easy way to create this file is to use `pg_dump` to dump out the table declarations and data needed

to set the scene, then add the problem query. You are encouraged to minimize the size of your example, but this is not absolutely necessary. If the bug is reproducible, we will find it either way.

If your application uses some other client interface, such as PHP, then please try to isolate the offending queries.

- The output you got. If there is an error message, show it. If the program terminates with an operating system error, say which. If nothing at all happens, say so. Even if the result of your test case is a program crash or otherwise obvious it might not happen on our platform. The easiest thing is to copy the output from the terminal, if possible.

Note

If you are reporting an error message, please obtain the most verbose form of the message. In `psql`, say `\set VERBOSITY verbose` beforehand. If you are extracting the message from the server log, set the run-time parameter `log_error_verbosity` to `verbose` so that all details are logged.

Note

In case of fatal errors, the error message reported by the client might not contain all the information available. Please also look at the log output of the database server.

- The output you expect is very important to state. Please provide the expected output, if applicable.
- Any command line options and other start-up options, including any relevant environment variables or configuration files that you changed from the default.
- Anything you did at all differently from the installation instructions.
- The Postgres Pro version. You can run the command `SELECT pgpro_version();` to find out the version of the server you are connected to. Most executable programs also support a `--version` option; at least `postgres --version` and `psql --version` should work.
- Platform information. This includes the kernel name and version, C library, processor, memory information, and so on.

5.3. Where to Report Bugs

In general, send bug reports to our support email address at `<bugs@postgrespro.ru>`. You are requested to use a descriptive subject for your email message, perhaps parts of the error message.

Do not send bug reports specific to Postgres Pro to the PostgreSQL support email address, as Postgres Pro is not supported by the PostgreSQL community. But you can send reports to `<pgsql-bugs@lists.postgresql.org>` for any bugs related to PostgreSQL.

Even if your bug is not specific to Postgres Pro, do not send bug reports to any of the user mailing lists, such as `<pgsql-sql@lists.postgresql.org>` or `<pgsql-general@lists.postgresql.org>`. These mailing lists are for answering user questions, and their subscribers normally do not wish to receive bug reports. More importantly, they are unlikely to fix them.

Also, please do *not* send reports to the developers' mailing list `<pgsql-hackers@lists.postgresql.org>`. This list is for discussing the development of PostgreSQL, and it would be nice if the community could keep the bug reports separate. The community might choose to take up a discussion about your bug report on `pgsql-hackers`, if the PostgreSQL-related problem needs more review.

Part I. Tutorial

Welcome to the Postgres Pro Tutorial. The tutorial is intended to give an introduction to Postgres Pro, relational database concepts, and the SQL language. We assume some general knowledge about how to use computers and no particular Unix or programming experience is required. This tutorial is intended to provide hands-on experience with important aspects of the Postgres Pro system. It makes no attempt to be a comprehensive treatment of the topics it covers.

After you have successfully completed this tutorial you will want to read the [Part II](#) section to gain a better understanding of the SQL language, or [Part IV](#) for information about developing applications with Postgres Pro. When learning SQL, you can use the demo database described in [Appendix K](#). Those who set up and manage their own server should also read [Part III](#).

Chapter 1. Getting Started

1.1. Quickstart Guide

If your site administrator has not set things up in the default way, you might have some more work to do. For example, if the database server machine is a remote machine, you will need to set the `PGHOST` environment variable to the name of the database server machine. The environment variable `PGPORT` might also have to be set. The bottom line is this: if you try to start an application program and it complains that it cannot connect to the database, you should consult your site administrator or, if that is you, the documentation to make sure that your environment is properly set up. If you did not understand the preceding paragraph then read the next section.

Postgres Pro Shardman is composed of several software components:

- PostgreSQL 17 DBMS with a set of patches.
- Postgres Pro Shardman extension.
- Management tools and services, including built-in shard manager to provide high availability.
- BiHA extension.

Postgres Pro Shardman stores its configuration in an etcd cluster. Therefore, we can use an existing etcd cluster, or we can deploy a simple one-node etcd cluster.

The [shardmand](#) daemon monitors the cluster configuration and manages BiHA clusters, which are used to guarantee high availability of all shards. The common Postgres Pro Shardman configuration (shardmand) is stored in an etcd cluster.

Currently Postgres Pro Shardman packages are available for

- Ubuntu 20.04/22.04
- Debian 10/11/12
- Red Hat Enterprise Linux 7/8/9
- Red OS 7.3/8
- Alt 9/10/10SP
- Astra Linux 1.7/1.8/s390

1.1.1. Cluster Configuration

Assume that we have three nodes for deploying Postgres Pro Shardman. Let's make the first one for the etcd one-node cluster and the other two nodes for the Postgres Pro Shardman two-node cluster.

Let's suppose that we have the following node names and IP addresses:

```
192.0.1.1 etcd - etcd one-node cluster
192.0.1.20 sdm01 - Shardman node1
192.0.1.21 sdm02 - Shardman node2
```

Each node has 4Gb RAM, 20GB HDD, 2CPU and Ubuntu 22.04 installed.

1.1.2. Preparation

1.1.2.1. Add host names to `/etc/hosts`

This step must be performed on all nodes.

```
sudo /bin/sh -c 'cat << EOF >> /etc/hosts
192.0.1.1 etcd
192.0.1.20 sdm01
```

```
192.0.1.21 sdm02
EOF'
```

1.1.2.2. Time Synchronization

This step must be performed on all nodes.

Deploy and start chrony daemon on all hosts.

```
sudo apt install -y chrony
```

By default, chrony gets the time from available servers on internet or the local time server. You can check available time servers as follows:

```
chronyc sources
MS Name/IP address          Stratum Poll Reach LastRx Last sample
=====
^? 192.0.1.100              1    6    7    1   -98us[ -98us] +/-  11ms
^* time.cloudflare.com      3    6    7    1   +139us[ +163us] +/-  11ms
^+ tms04.deltatelesystems.ru 1    6    7    1   -381us[ -357us] +/-  17ms
```

It is desirable to synchronize time with your server or the local server for the cluster. To do this, make changes similar to the following to chrony configuration:

```
cat /etc/chrony/chrony.conf

server 192.0.1.100 iburst
keyfile /etc/chrony.keys
driftfile /var/lib/chrony/chrony.drift
log tracking measurements statistics
logdir /var/log/chrony
```

```
systemctl restart chrony
```

Check that chrony is connected to the appropriate server.

```
chronyc sources
MS Name/IP address          Stratum Poll Reach LastRx Last sample
=====
^? 192.0.1.100              8    6   17   37   +14us[ +70us] +/- 161us
chronyc tracking
Reference ID      : 0A80000C (ntp.local)
Stratum          : 9
Ref time (UTC)   : Wed Nov 15 11:58:52 2023
System time      : 0.000000004 seconds slow of NTP time
Last offset      : -0.000056968 seconds
RMS offset       : 0.000056968 seconds
Frequency        : 10.252 ppm fast
Residual freq    : -2.401 ppm
Skew             : 364.419 ppm
Root delay       : 0.000455358 seconds
Root dispersion  : 0.010503666 seconds
Update interval  : 2.1 seconds
Leap status      : Normal
```

1.1.3. Deploy an etcd One-Node Cluster

Note also a [Deploy a Multi-Node etcd cluster](#) section.

Install the following packages:

```
sudo apt install -y vim curl
```

- Run

- Open the file `pgpro-repo-add.sh` and specify the repository password in the `PASSWORD` variable.
- Run `sudo pgpro-repo-add.sh`.

```
sudo apt install -y etcd-sdm
```

```
sudo vim /etc/default/etcd-sdm
ETCD_NAME=etcd
ETCD_LISTEN_CLIENT_URLS=http://0.0.0.0:2379
ETCD_ADVERTISE_CLIENT_URLS=http://192.0.1.1:2379
ETCD_MAX_SNAPSHOTS=5
ETCD_MAX_WALS=5
ETCD_AUTO_COMPACTION_MODE=periodic
ETCD_AUTO_COMPACTION_RETENTION=5m
ETCD_QUOTA_BACKEND_BYTES=6442450944
ETCD_DATA_DIR=/var/lib/etcd-sdm/sdm-17
```

Clear the etcd data directory:

Restart the etcd-sdm service:

For your user, add `/opt/pgpro/sdm-17/bin` to the `PATH` environment variable:

Check that etcd is properly configured:

4

etcd one-node cluster is properly configured and ready to serve requests.

To prevent bloat when etcd is intensively used, add a defragmentation command to cron:

```
sudo sh -c '{ crontab -l; echo "@hourly /opt/pgpro/sdm-17/bin/etcdctl defrag"; }
| crontab'
```

1.1.4. Deploy Postgres Pro Shardman Nodes

Let's add a Postgres Pro Shardman repository on each node:

- Set a user and password as in [Section 1.1.3](#).
- Run

```
curl -fsSL -u "<user>:<password>" https://repo.postgrespro.ru/sdm/sdm-17/keys/pgpro-
repo-add.sh > pgpro-repo-add.sh | bash
chmod +x pgpro-repo-add.sh
```

- Open the file `pgpro-repo-add.sh` and specify the repository password in the `PASSWORD` variable.
- Run `pgpro-repo-add.sh`.

Next step is installation of packages (on each node):

```
sudo apt update
sudo apt install -y postgrespro-sdm-17-server postgrespro-sdm-17-client postgrespro-
sdm-17-contrib postgrespro-sdm-17-libs pg-probackup-sdm-17 shardman-services shardman-
tools
```

Suppose we have chosen a default cluster name of `cluster0`. The next step is to put Postgres Pro Shardman environment vars into the `/etc/shardman` directory (on each node):

```
sudo sh -c 'cat << EOF > /etc/shardman/shardmand-cluster0.env
SDM_CLUSTER_NAME=cluster0
SDM_LOG_LEVEL=info
SDM_STORE_ENDPOINTS=http://etcd:2379
EOF'
```

The file and directory are created with `sudo`, but later `shardmanctl` does not use `sudo`, thus cannot access the file with the environment variables. To access it, either add the variables to the system with `export`, or grant user with access rights to the file and the directory.

For your user, add `/opt/pgpro/sdm-17/bin` to the `PATH` environment variable and export the `SDM_STORE_ENDPOINTS` variable:

```
echo "export PATH=$PATH:/opt/pgpro/sdm-17/bin" >> ~/.bashrc
echo "export SDM_STORE_ENDPOINTS=http://etcd:2379" >> ~/.bashrc
source ~/.bashrc
```

Let's generate a sample configuration with the Postgres Pro Shardman utilities (only on one node).

```
shardmanctl config generate > spec.json
```

In this step, you can make some changes to the cluster specification (configuration), i.e., change the password or PostgreSQL `shared_buffers` parameter and so on.

1.1.5. Initialize the Postgres Pro Shardman Cluster

Now we have some final steps. First, let's initialize the cluster configuration in etcd (only on one [any] node).

```
shardmanctl init -f spec.json
```

The expected output is:

```
2024-10-22T11:34:05.713Z      DEBUG    Waiting for metadata lock...    {"caller":  
  "cmd/common.go:198"}  
2024-10-22T11:34:05.721Z      DEBUG    DataDir is not specified, setting to default /  
var/lib/pgpro/sdm-17/data {"caller": "cluster/cluster.go:860"}
```

Enable and start the shardmand service (on each node):

```
sudo systemctl enable --now shardmand@cluster0  
sudo systemctl status shardmand@cluster0  
  
● shardmand@cluster0.service - deployment daemon for shardman  
   Loaded: loaded (/lib/systemd/system/shardmand@.service; enabled; vendor preset:  
enabled)  
   Active: active (running) since Tue 2023-04-18 12:28:18 UTC; 2min 13s ago  
     Docs: https://github.com/postgrespro/shardman  
    Main PID: 618 (shardmand)  
      Tasks: 10 (limit: 4571)  
     Memory: 32.0M  
        CPU: 422ms  
    CGroup: /system.slice/system-shardmand.slice/shardmand@cluster0.service  
            └─618 /opt/pgpro/sdm-17/bin/shardmand --cluster-name cluster0 --system-  
bus --user postgres
```

1.1.6. Add Nodes to the Postgres Pro Shardman Cluster

In this step we assume that all previous steps were executed successfully: etcd cluster is working properly, the time on all hosts is synchronized, and the daemon is launched on `sdm01` and `sdm02`. The final step should be executed with `shardmanctl` command as follows:

```
shardmanctl nodes add -n sdm01,sdm02 \  
    --cluster-name cluster0 \  
    --log-level debug
```

The expected output should be:

```
2023-04-18T12:43:11.300Z DEBUG cmd/common.go:100 Waiting for metadata lock...  
2023-04-18T12:43:11.306Z INFO cluster/store.go:277 Checking if shardmand on all nodes  
have applied current cluster configuration  
# Waiting for shardmand on node sdm01 to apply current configuration: success 0.000s  
# Waiting for shardmand on node sdm02 to apply current configuration: success 0.000s  
2023-04-18T12:43:11.307Z INFO add/case.go:112 Initting Stolon instances...  
2023-04-18T12:43:11.312Z INFO add/case.go:170 Waiting for Stolon daemons to start...  
make sure shardmand daemons are running on the nodes  
# Waiting for Stolon daemons of rg clover-1-sdm01: success 31.012s  
# Waiting for Stolon daemons of rg clover-1-sdm02: success 0.012s  
2023-04-18T12:43:42.336Z INFO add/case.go:187 Adding repgroups...  
# waiting rg 1 config apply: done 7.014s  
2023-04-18T12:43:49.444Z DEBUG broadcaster/worker.go:33 start broadcaster worker for  
repgroup id=1  
2023-04-18T12:43:49.453Z DEBUG broadcaster/worker.go:51 repgroup 1 connect  
established  
2023-04-18T12:43:49.453Z DEBUG commands/addrepgroup.go:575 waiting for extension  
lock...  
2023-04-18T12:43:49.453Z DEBUG commands/addrepgroup.go:137 Loading schema into  
replication group rg 1  
...  
2023-04-18T12:44:25.665Z DEBUG rebalance/service.go:528 wait all tasks finish
```

```
2023-04-18T12:44:25.666Z DEBUG broadcaster/worker.go:75 finish broadcaster worker for
repgroup id=1
2023-04-18T12:44:25.666Z DEBUG broadcaster/worker.go:75 finish broadcaster worker for
repgroup id=2
2023-04-18T12:44:25.666Z INFO add/case.go:221 Successfully added nodes sdm01, sdm02
to the cluster
```

The “Successfully added nodes sdm01, sdm02 to the cluster” message means that everything is fine and nodes sdm01 and sdm02 are working properly.

1.1.7. Check the Postgres Pro Shardman Cluster Status

Let's check the status of the cluster nodes.

```
shardmanctl status
```

== STORE STATUS ==			
STATUS	NODE	MESSAGE	REPLICATION GROUP
Warning		Store has only one member, consider deploying store cluster	
== METADATA STATUS ==			
STATUS	NODE	MESSAGE	REPLICATION GROUP
OK		Metadata is OK	
== SHARDMAND STATUS ==			
STATUS	NODE	MESSAGE	REPLICATION GROUP
OK	sdm01	shardmand status on node sdm01 is OK	
OK	sdm02	shardmand status on node sdm02 is OK	
== REPLICATION GROUP STATUS ==			
STATUS	NODE	MESSAGE	REPLICATION GROUP
Error		Replication group shard-1 is OK	shard-1

Getting Started

Error	Replication group shard-2 is OK	shard-2
-------	---------------------------------	---------

== MASTER STATUS ==

STATUS	NODE	MESSAGE	REPLICATION GROUP
OK	sdm01:5432	Replication group shard-1 master is running on sdm01:5432	shard-1
OK	sdm02:5432	Replication group shard-2 master is running on sdm02:5432	shard-2

== DICTIONARY STATUS ==

STATUS	NODE	MESSAGE	REPLICATION GROUP
OK		Replication group shard-2 dictionary is OK	shard-2
OK		Replication group shard-1 dictionary is OK	shard-1

== KEEPER STATUS ==

STATUS	NODE	MESSAGE	REPLICATION GROUP
OK	sdm01:5432	Keeper keeper_1 is OK	shard-1
OK	sdm02:5432	Keeper keeper_1 is OK	shard-2

1.1.8. Connect to the Postgres Pro Shardman Cluster

To connect to the cluster we should get the cluster connection string on any cluster node (sdm01 or sdm02):

```
shardmanctl getconnstr
```

```
dbname=postgres host=sdm01,sdm02 password=!!!CHANGE_ME!!! port=5432,5432  
user=postgres
```

And then let's try to connect:


```
psql -d 'dbname=postgres host=sdm01,sdm02 password=!!!CHANGE_ME!!! port=5432,5432
user=postgres'

psql (14.7)
Type "help" for help.

postgres=#
```

1.1.9. Example: Deploy a Multi-Node etcd Cluster

The process is described for the following servers:

```
192.0.1.1 etcd1
192.0.1.2 etcd2
192.0.1.3 etcd3
```

Install the needed packages on each server:

```
sudo apt install -y vim curl
```

To connect the repository, on each server, run:

```
sudo curl -fsSL https://repo.postgrespro.ru/sdm/sdm-17/keys/pgpro-repo-add.sh | bash
```

Install etcd-sdm packages on each server:

```
sudo apt install -y etcd-sdm
```

For each server, edit the file that lists environment variables, replacing placeholders in angle brackets with specific values:

```
sudo vim /etc/default/etcd-sdm
ETCD_NAME=<hostname>
ETCD_LISTEN_PEER_URLS=http://0.0.0.0:2380
ETCD_LISTEN_CLIENT_URLS=http://0.0.0.0:2379
ETCD_ADVERTISE_CLIENT_URLS=http://<host ip address>:2379
ETCD_INITIAL_ADVERTISE_PEER_URLS=http://<host ip address>:2380
ETCD_INITIAL_CLUSTER_TOKEN=etcd-cluster-1
ETCD_INITIAL_CLUSTER_STATE=new
ETCD_MAX_SNAPSHOTS=5
ETCD_MAX_WALS=5
ETCD_AUTO_COMPACTION_MODE=periodic
ETCD_AUTO_COMPACTION_RETENTION=5m
ETCD_QUOTA_BACKEND_BYTES=6442450944
ETCD_DATA_DIR=/var/lib/etcd-sdm/sdm-17
ETCD_INITIAL_CLUSTER=etcd1=http://<ip etcd1>:2380,etcd2=http://<ip
etcd2>:2380,etcd3=http://<ip etcd3>:2380
```

This file will be loaded at etcd start with its own start settings on each server.

Clear the etcd data directory:

```
sudo rm -rf /var/lib/etcd-sdm/sdm-17/*
```

Restart the etcd-sdm service on each server:

```
sudo systemctl restart etcd-sdm
```

For your user, add /opt/pgpro/sdm-17/bin to the PATH environment variable:

```
echo "export PATH=$PATH:/opt/pgpro/sdm-17/bin" >> ~/.bashrc
```

Check that etcd is properly configured:

```

+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|          ID          | STATUS  | NAME    |          PEER ADDRS          |          CLIENT          |
| ADDRS              | IS LEARNER |         |                               |                           |
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
| 318be6342e6d9ac | started | etcd1   | http://192.0.1.1:2380      |                           |
| http://192.0.1.1:2379 |         | false   |                               |                           |
| 9e49480544aedb89 | started | etcd2   | http://192.0.1.2:2380      |                           |
| http://192.0.1.2:2379 |         | false   |                               |                           |
| bb3772bfa22482d7 | started | etcd3   | http://192.0.1.3:2380      |                           |
| http://192.0.1.3.4:2379 |         | false   |                               |                           |
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
$ etcdctl --endpoints=http://192.0.1.1:2380,http://192.0.1.2:2380,http://192.0.1.3:2380
endpoint status health -w table
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|          ENDPOINT          |          ID          | VERSION | DB SIZE | IS LEADER | IS
| LEARNER | RAFT TERM | RAFT INDEX | RAFT APPLIED INDEX | ERRORS |
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
| http://192.0.1.1:2380      | 318be6342e6d9ac | 3.5.9 | 5.7 MB | true |
| false | 13 | 425686 | 425686 | |
| http://192.0.1.2:2380      | 9e49480544aedb89 | 3.5.9 | 5.7 MB | false |
| false | 13 | 425686 | 425686 | |
| http://192.0.1.3:2380      | bb3772bfa22482d7 | 3.5.9 | 5.7 MB | false |
| false | 13 | 425686 | 425686 | |
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+

```

To prevent bloat when etcd is intensively used, add a defragmentation command to cron:

The final endpoints string of the etcd cluster:

It should be specified in `/etc/shardman` configuration file and as a `--store-endpoints` parameter of `shardmanctl`.

1.2. When to use

Postgres Pro Shardman provides horizontal scalability with a view and consistency of a single database. Applications can use every node to access the distributed database and operate mostly the same way as with a single PostgreSQL instance. Still internally it is a distributed system that imposes certain rules on designing schema and writing queries. The main direction of adoption is to localize the data and the computations.

10

- The working set of data does not fit in RAM of one server. Distributed systems can have much bigger total size of RAM.
- Maintenance operations such as vacuum take too long. Shardman utilizes partitioned tables under the hood. Maintenance operations can be parallelized by nodes and partitions of tables.
- Number of read sessions is too large for one instance of PostgreSQL. Shardman allows to distribute read sessions across the cluster and handle internal connections very efficiently with multiplexing transport.
- Intensive write operations. Distributed systems can have much bigger total number of disk IOPS.
- CPU intensive queries. Shardman allows to distribute calculations by nodes and reduce execution time for complex queries.

When Postgres Pro Shardman is not appropriate:

- Vertical scaling is economically and technically possible.
- Data model and workload require a lot of cross-shard transactions.
- Complex analytics, in particular joins of sharded tables when conditions don't include the sharding key.
- Multi-DC/Multi-region deployments.

1.3. Architectural Fundamentals

Before we proceed, you should understand the basic Postgres Pro system architecture. Understanding how the parts of Postgres Pro interact will make this chapter somewhat clearer.

In database jargon, Postgres Pro uses a client/server model. A Postgres Pro session consists of the following cooperating processes (programs):

- A server process, which manages the database files, accepts connections to the database from client applications, and performs database actions on behalf of the clients. The database server program is called `postgres`.
- The user's client (frontend) application that wants to perform database operations. Client applications can be very diverse in nature: a client could be a text-oriented tool, a graphical application, a web server that accesses the database to display web pages, or a specialized database maintenance tool. Some client applications are supplied with the Postgres Pro distribution; most are developed by users.

As is typical of client/server applications, the client and the server can be on different hosts. In that case they communicate over a TCP/IP network connection. You should keep this in mind, because the files that can be accessed on a client machine might not be accessible (or might only be accessible using a different file name) on the database server machine.

The Postgres Pro server can handle multiple concurrent connections from clients. To achieve this it starts (“forks”) a new process for each connection. From that point on, the client and the new server process communicate without intervention by the original `postgres` process. Thus, the supervisor server process is always running, waiting for client connections, whereas client and associated server processes come and go. (All of this is of course invisible to the user. We only mention it here for completeness.)

1.4. Creating a Database

You can only create local DBs, as `postgres` is the only distributed DB.

The first test to see whether you can access the database server is to try to create a database. A running Postgres Pro server can manage many databases. Typically, a separate database is used for each project or for each user.

Possibly, your site administrator has already created a database for your use. In that case you can omit this step and skip ahead to the next section.

To create a new database, in this example named `mydb`, you use the following command:

```
$ createdb mydb
```

If this produces no response then this step was successful and you can skip over the remainder of this section.

If you see a message similar to:

```
createdb: command not found
```

then Postgres Pro was not installed properly. Either it was not installed at all or your shell's search path was not set to include it. Try calling the command with an absolute path instead:

```
$ /usr/local/pgsql/bin/createdb mydb
```

The path at your site might be different. Contact your site administrator or check the installation instructions to correct the situation.

Another response could be this:

```
createdb: error: connection to server on socket "/tmp/.s.PGSQL.5432" failed: No such file or directory
```

```
Is the server running locally and accepting connections on that socket?
```

This means that the server was not started, or it is not listening where `createdb` expects to contact it. Again, check the installation instructions or consult the administrator.

Another response could be this:

```
createdb: error: connection to server on socket "/tmp/.s.PGSQL.5432" failed: FATAL: role "joe" does not exist
```

where your own login name is mentioned. This will happen if the administrator has not created a Postgres Pro user account for you. (Postgres Pro user accounts are distinct from operating system user accounts.) If you are the administrator, see [Chapter 21](#) for help creating accounts. You will need to become the operating system user under which Postgres Pro was installed (usually `postgres`) to create the first user account. It could also be that you were assigned a Postgres Pro user name that is different from your operating system user name; in that case you need to use the `-U` switch or set the `PGUSER` environment variable to specify your Postgres Pro user name.

If you have a user account but it does not have the privileges required to create a database, you will see the following:

```
createdb: error: database creation failed: ERROR: permission denied to create database
```

Not every user has authorization to create new databases. If Postgres Pro refuses to create databases for you then the site administrator needs to grant you permission to create databases. Consult your site administrator if this occurs. If you installed Postgres Pro yourself then you should log in for the purposes of this tutorial under the user account that you started the server as.¹

You can also create databases with other names. Postgres Pro allows you to create any number of databases at a given site. Database names must have an alphabetic first character and are limited to 63 bytes in length. A convenient choice is to create a database with the same name as your current user name. Many tools assume that database name as the default, so it can save you some typing. To create that database, simply type:

```
$ createdb
```

If you do not want to use your database anymore you can remove it. For example, if you are the owner (creator) of the database `mydb`, you can destroy it using the following command:

```
$ dropdb mydb
```

¹ As an explanation for why this works: Postgres Pro user names are separate from operating system user accounts. When you connect to a database, you can choose what Postgres Pro user name to connect as; if you don't, it will default to the same name as your current operating system account. As it happens, there will always be a Postgres Pro user account that has the same name as the operating system user that started the server, and it also happens that that user always has permission to create databases. Instead of logging in as that user you can also specify the `-U` option everywhere to select a Postgres Pro user name to connect as.

(For this command, the database name does not default to the user account name. You always need to specify it.) This action physically removes all files associated with the database and cannot be undone, so this should only be done with a great deal of forethought.

More about `createdb` and `dropdb` can be found in [createdb](#) and [dropdb](#) respectively.

1.5. Accessing a Database

As Postgres Pro Shardman is a distributed system, it is sufficient to connect to one cluster node to connect to entire cluster. See [connection details](#). Once you have created a database, you can access it by:

- Running the Postgres Pro interactive terminal program, called *psql*, which allows you to interactively enter, edit, and execute SQL commands.
- Writing a custom application, using one of the several available language bindings. These possibilities are discussed further in [Part IV](#).

You probably want to start up *psql* to try the examples in this tutorial. It can be activated for the `mydb` database by typing the command:

```
$ psql mydb
```

If you do not supply the database name then it will default to your user account name. You already discovered this scheme in the previous section using `createdb`.

In *psql*, you will be greeted with the following message:

```
psql (17.5.1)
Type "help" for help.
```

```
mydb=>
```

The last line could also be:

```
mydb=#
```

That would mean you are a database superuser, which is most likely the case if you installed the Postgres Pro instance yourself. Being a superuser means that you are not subject to access controls. For the purposes of this tutorial that is not important.

If you encounter problems starting *psql* then go back to the previous section. The diagnostics of `createdb` and *psql* are similar, and if the former worked the latter should work as well.

The last line printed out by *psql* is the prompt, and it indicates that *psql* is listening to you and that you can type SQL queries into a work space maintained by *psql*. Try out these commands:

```
mydb=> SELECT pgpro_version();
                                version
-----
Postgres Pro Shardman 17.5.1 on x86_64-pc-linux-gnu, compiled by gcc (Ubuntu
11.4.0-1ubuntu1~22.04) 11.4.0, 64-bit
(1 row)

mydb=> SELECT current_date;
      date
-----
2016-01-07
(1 row)

mydb=> SELECT 2 + 2;
?column?
-----
4
```

```
(1 row)
```

The `psql` program has a number of internal commands that are not SQL commands. They begin with the backslash character, “\”. For example, you can get help on the syntax of various Postgres Pro SQL commands by typing:

```
mydb=> \h
```

To get out of `psql`, type:

```
mydb=> \q
```

and `psql` will quit and return you to your command shell. (For more internal commands, type `\?` at the `psql` prompt.) The full capabilities of `psql` are documented in [psql](#). In this tutorial we will not use these features explicitly, but you can use them yourself when it is helpful.

Note that to connect to a Postgres Pro Shardman cluster and successfully work with it, it is sufficient to connect to one cluster node.

1.5.1. `psql/libpq`

To connect to a Postgres Pro Shardman cluster, first get the [connection string](#) with `psql` or `libpq`. For more information on `libpq`, refer to [libpq — C Library](#).

The value of `target_session_attrs` must be set to `read-write`. Only connections that allow read/write transactions are acceptable. If the connection to a cluster node is a success, the request “SHOW transaction_read_only;” is sent. If it returns `on`, the connection is closed. If several servers are specified in the connection string, other servers will be iterated through, the same way as with the failed connection attempt. The `target_session_attrs` parameter allows you to specify both masters and replicas of the Postgres Pro Shardman cluster.

The following examples illustrate the connection:

```
psql -d "dbname=postgres host=node3,node4,node2,node1 port=5432,5432,5432,5432
user=username password=password target_session_attrs=read-write"
```

```
psql postgres://username:password@node1:5432,node2:5432,node3:5432,node4:5432/postgres?
target_session_attrs=read-write
```

1.5.2. Python

Connection to a Postgres Pro Shardman cluster using the `psycopg2` library looks like this:

```
import psycopg2
from psycopg2 import pool
```

```
pool = psycopg2.pool.SimpleConnectionPool(
    min_size=1,
    max_size=5,
    user="pguser",
    password="*****",
    host="node1,node2,node3,node4",
    port="5432,5432,5432,5432",
    database="postgres",
    target_session_attrs="read-write")
```

```
connection = pool.getconn()
```

A connection pool with the following parameters is created: the minimum and maximum number of connections `min_size=1` and `max_size=5`. Then a specific connection to the cluster is selected, the user login and password are specified, as well as the list of nodes and TCP ports, database and connection parameters (see [Section 1.5.1](#) for more information).

1.5.3. Java

Connection to a Postgres Pro Shardman cluster using JDBC looks like this:

```
String url = "jdbc:postgresql://node1:5432,node2:5432,node3:5432,node4:5432/postgres?loadBalanceHosts=true&targetServerType=primary";
Properties props = new Properties();
```

```
props.setProperty("user", "postgres");
props.setProperty("password", "*****");
```

```
Connection conn = DriverManager.getConnection(url, props);
```

`url` contains the connection string, where all the available shard masters are listed. If no additional connection parameters of the JDBC driver are specified, connection to the cluster is performed through the first node available for connection. This is not always convenient. Therefore, connection string settings are added that allow using different cluster shards for different connections.

`loadBalanceHosts=true` allows iterating through nodes connecting to one of them, and `targetServerType=primary` indicates a need to only choose masters, then replicas can be added to the connection string.

1.5.4. Go

Ways to connect to a Postgres Pro Shardman cluster for Go are pretty much the same as those accepted in Java or Python. You need to specify lists of nodes, their TCP ports, as well as connection parameters and choose a suitable driver.

One of these drivers for Go is [pgx](#) version 4 or 5.

The following is an example of a connection string and creation of a pool for connecting to a cluster:

```
dbURL := "postgres://username:password@node1:5432,node2:5432,node3:5432,node4:5432/postgres?target_session_attrs=read-write"
dbPool, err := pgxpool.New(context.Background(), dbURL)
```

Also pay attention to the [description of the `target_session_attrs` parameter](#).

Chapter 2. The SQL Language

2.1. Introduction

This chapter provides an overview of how to use SQL to perform simple operations. This tutorial is only intended to give you an introduction and is in no way a complete tutorial on SQL. Numerous books have been written on SQL, including [melt93](#) and [date97](#). You should be aware that some Postgres Pro language features are extensions to the standard.

In the examples that follow, we assume that you have created a database named `mydb`, as described in the previous chapter, and have been able to start `psql`.

2.2. Concepts

Postgres Pro is a *relational database management system* (RDBMS). That means it is a system for managing data stored in *relations*. Relation is essentially a mathematical term for *table*. The notion of storing data in tables is so commonplace today that it might seem inherently obvious, but there are a number of other ways of organizing databases. Files and directories on Unix-like operating systems form an example of a hierarchical database. A more modern development is the object-oriented database.

Each table is a named collection of *rows*. Each row of a given table has the same set of named *columns*, and each column is of a specific data type. Whereas columns have a fixed order in each row, it is important to remember that SQL does not guarantee the order of the rows within the table in any way (although they can be explicitly sorted for display).

Tables are grouped into databases, and a collection of databases managed by a single Postgres Pro server instance constitutes a database *cluster*.

2.3. Creating a New Table

You can create a new table by specifying the table name, along with all column names and their types.

2.3.1. Create Sharded Tables

Let's try to create a sharded table and check if everything is working properly.

```
postgres=# create table x(id int primary key, t text) with
(distributed_by='id',num_parts=2);
CREATE TABLE
```

```
postgres=# \d
```

```
                List of relations
Schema | Name      | Type              | Owner
-----+-----+-----+-----
public | x         | partitioned table | postgres
public | x_0       | table             | postgres
public | x_1_fdw   | foreign table     | postgres
(3 rows)
```

```
postgres=# \d x_0
```

```
Table "public.x_0"
Column | Type   | Collation | Nullable | Default
-----+-----+-----+-----+-----
id     | integer |           | not null |
t      | text   |           |          |
Partition of: x FOR VALUES WITH (modulus 2, remainder 0)
Indexes:
```



```
"x_0_pkey" PRIMARY KEY, btree (id)

postgres=# \d x_1_fdw
               Foreign table "public.x_1_fdw"
  Column | Type   | Collation | Nullable | Default | FDW options
-----+-----+-----+-----+-----+-----
  id     | integer |           | not null |         |
  t      | text    |           |         |         |
Partition of: x FOR VALUES WITH (modulus 2, remainder 1)
Server: shardman_rg_2
FDW options: (table_name 'x_1')

postgres=# insert into x values (1,'t'), (2,'t'), (3,'t');
INSERT 0 3

postgres=# select * from x_0;
 id | t
----+---
  1 | t
  2 | t
(2 rows)

postgres=# select * from x_1_fdw;
 id | t
----+---
  3 | t
(1 row)
```

Everything works as expected.

2.3.2. Create Regular Tables

For regular tables:

```
CREATE TABLE weather (
    city          varchar(80),
    temp_lo       int,          -- low temperature
    temp_hi       int,          -- high temperature
    prcp          real,         -- precipitation
    date          date
);
```

You can enter this into `psql` with the line breaks. `psql` will recognize that the command is not terminated until the semicolon.

White space (i.e., spaces, tabs, and newlines) can be used freely in SQL commands. That means you can type the command aligned differently than above, or even all on one line. Two dashes (“--”) introduce comments. Whatever follows them is ignored up to the end of the line. SQL is case-insensitive about key words and identifiers, except when identifiers are double-quoted to preserve the case (not done above).

`varchar(80)` specifies a data type that can store arbitrary character strings up to 80 characters in length. `int` is the normal integer type. `real` is a type for storing single precision floating-point numbers. `date` should be self-explanatory. (Yes, the column of type `date` is also named `date`. This might be convenient or confusing — you choose.)

Postgres Pro supports the standard SQL types `int`, `smallint`, `real`, `double precision`, `char(N)`, `varchar(N)`, `date`, `time`, `timestamp`, and `interval`, as well as other types of general utility and a rich set of geometric types. Postgres Pro can be customized with an arbitrary number of user-defined data types. Consequently, type names are not key words in the syntax, except where required to support special cases in the SQL standard.

The second example will store cities and their associated geographical location:

```
CREATE TABLE cities (  
    name          varchar(80),  
    location      point  
);
```

The `point` type is an example of a Postgres Pro-specific data type.

Finally, it should be mentioned that if you don't need a table any longer or want to recreate it differently you can remove it using the following command:

```
DROP TABLE tablename;
```

2.4. Populating a Table With Rows

The `INSERT` statement is used to populate a table with rows:

```
INSERT INTO weather VALUES ('San Francisco', 46, 50, 0.25, '1994-11-27');
```

Note that all data types use rather obvious input formats. Constants that are not simple numeric values usually must be surrounded by single quotes (`'`), as in the example. The `date` type is actually quite flexible in what it accepts, but for this tutorial we will stick to the unambiguous format shown here.

The `point` type requires a coordinate pair as input, as shown here:

```
INSERT INTO cities VALUES ('San Francisco', '(-194.0, 53.0)');
```

The syntax used so far requires you to remember the order of the columns. An alternative syntax allows you to list the columns explicitly:

```
INSERT INTO weather (city, temp_lo, temp_hi, prcp, date)  
VALUES ('San Francisco', 43, 57, 0.0, '1994-11-29');
```

You can list the columns in a different order if you wish or even omit some columns, e.g., if the precipitation is unknown:

```
INSERT INTO weather (date, city, temp_hi, temp_lo)  
VALUES ('1994-11-29', 'Hayward', 54, 37);
```

Many developers consider explicitly listing the columns better style than relying on the order implicitly.

Please enter all the commands shown above so you have some data to work with in the following sections.

You could also have used `COPY` to load large amounts of data from flat-text files. This is usually faster because the `COPY` command is optimized for this application while allowing less flexibility than `INSERT`. An example would be:

```
COPY weather FROM '/home/user/weather.txt';
```

where the file name for the source file must be available on the machine running the backend process, not the client, since the backend process reads the file directly. You can read more about the `COPY` command in [COPY](#).

2.5. Querying a Table

To retrieve data from a table, the table is *queried*. An SQL `SELECT` statement is used to do this. The statement is divided into a select list (the part that lists the columns to be returned), a table list (the part that lists the tables from which to retrieve the data), and an optional qualification (the part that specifies any restrictions). For example, to retrieve all the rows of table `weather`, type:

```
SELECT * FROM weather;
```

Here `*` is a shorthand for “all columns”.¹ So the same result would be had with:

¹ While `SELECT *` is useful for off-the-cuff queries, it is widely considered bad style in production code, since adding a column to the table would change the results.

```
SELECT city, temp_lo, temp_hi, prcp, date FROM weather;
```

The output should be:

city	temp_lo	temp_hi	prcp	date
San Francisco	46	50	0.25	1994-11-27
San Francisco	43	57	0	1994-11-29
Hayward	37	54		1994-11-29

(3 rows)

You can write expressions, not just simple column references, in the select list. For example, you can do:

```
SELECT city, (temp_hi+temp_lo)/2 AS temp_avg, date FROM weather;
```

This should give:

city	temp_avg	date
San Francisco	48	1994-11-27
San Francisco	50	1994-11-29
Hayward	45	1994-11-29

(3 rows)

Notice how the `AS` clause is used to relabel the output column. (The `AS` clause is optional.)

A query can be “qualified” by adding a `WHERE` clause that specifies which rows are wanted. The `WHERE` clause contains a Boolean (truth value) expression, and only rows for which the Boolean expression is true are returned. The usual Boolean operators (`AND`, `OR`, and `NOT`) are allowed in the qualification. For example, the following retrieves the weather of San Francisco on rainy days:

```
SELECT * FROM weather
WHERE city = 'San Francisco' AND prcp > 0.0;
```

Result:

city	temp_lo	temp_hi	prcp	date
San Francisco	46	50	0.25	1994-11-27

(1 row)

You can request that the results of a query be returned in sorted order:

```
SELECT * FROM weather
ORDER BY city;
```

city	temp_lo	temp_hi	prcp	date
Hayward	37	54		1994-11-29
San Francisco	43	57	0	1994-11-29
San Francisco	46	50	0.25	1994-11-27

In this example, the sort order isn't fully specified, and so you might get the San Francisco rows in either order. But you'd always get the results shown above if you do:

```
SELECT * FROM weather
ORDER BY city, temp_lo;
```

You can request that duplicate rows be removed from the result of a query:

```
SELECT DISTINCT city
FROM weather;
```

city
Hayward
San Francisco

(2 rows)

Here again, the result row ordering might vary. You can ensure consistent results by using `DISTINCT` and `ORDER BY` together:²

```
SELECT DISTINCT city
FROM weather
ORDER BY city;
```

2.6. Joins Between Tables

Thus far, our queries have only accessed one table at a time. Queries can access multiple tables at once, or access the same table in such a way that multiple rows of the table are being processed at the same time. Queries that access multiple tables (or multiple instances of the same table) at one time are called *join* queries. They combine rows from one table with rows from a second table, with an expression specifying which rows are to be paired. For example, to return all the weather records together with the location of the associated city, the database needs to compare the `city` column of each row of the `weather` table with the `name` column of all rows in the `cities` table, and select the pairs of rows where these values match.³ This would be accomplished by the following query:

```
SELECT * FROM weather JOIN cities ON city = name;
```

city	temp_lo	temp_hi	prcp	date	name	location
San Francisco	46	50	0.25	1994-11-27	San Francisco	(-194,53)
San Francisco	43	57	0	1994-11-29	San Francisco	(-194,53)

(2 rows)

Observe two things about the result set:

- There is no result row for the city of Hayward. This is because there is no matching entry in the `cities` table for Hayward, so the join ignores the unmatched rows in the `weather` table. We will see shortly how this can be fixed.
- There are two columns containing the city name. This is correct because the lists of columns from the `weather` and `cities` tables are concatenated. In practice this is undesirable, though, so you will probably want to list the output columns explicitly rather than using `*`:

```
SELECT city, temp_lo, temp_hi, prcp, date, location
FROM weather JOIN cities ON city = name;
```

Since the columns all had different names, the parser automatically found which table they belong to. If there were duplicate column names in the two tables you'd need to *qualify* the column names to show which one you meant, as in:

```
SELECT weather.city, weather.temp_lo, weather.temp_hi,
       weather.prcp, weather.date, cities.location
FROM weather JOIN cities ON weather.city = cities.name;
```

It is widely considered good style to qualify all column names in a join query, so that the query won't fail if a duplicate column name is later added to one of the tables.

Join queries of the kind seen thus far can also be written in this form:

```
SELECT *
FROM weather, cities
WHERE city = name;
```

This syntax pre-dates the `JOIN/ON` syntax, which was introduced in SQL-92. The tables are simply listed in the `FROM` clause, and the comparison expression is added to the `WHERE` clause. The results from this older

² In some database systems, including older versions of Postgres Pro, the implementation of `DISTINCT` automatically orders the rows and so `ORDER BY` is unnecessary. But this is not required by the SQL standard, and current Postgres Pro does not guarantee that `DISTINCT` causes the rows to be ordered.

³ This is only a conceptual model. The join is usually performed in a more efficient manner than actually comparing each possible pair of rows, but this is invisible to the user.

implicit syntax and the newer explicit `JOIN/ON` syntax are identical. But for a reader of the query, the explicit syntax makes its meaning easier to understand: The join condition is introduced by its own key word whereas previously the condition was mixed into the `WHERE` clause together with other conditions.

Now we will figure out how we can get the Hayward records back in. What we want the query to do is to scan the `weather` table and for each row to find the matching `cities` row(s). If no matching row is found we want some “empty values” to be substituted for the `cities` table's columns. This kind of query is called an *outer join*. (The joins we have seen so far are *inner joins*.) The command looks like this:

```
SELECT *
FROM weather LEFT OUTER JOIN cities ON weather.city = cities.name;
```

city	temp_lo	temp_hi	prcp	date	name	location
Hayward	37	54		1994-11-29		
San Francisco	46	50	0.25	1994-11-27	San Francisco	(-194,53)
San Francisco	43	57	0	1994-11-29	San Francisco	(-194,53)

(3 rows)

This query is called a *left outer join* because the table mentioned on the left of the join operator will have each of its rows in the output at least once, whereas the table on the right will only have those rows output that match some row of the left table. When outputting a left-table row for which there is no right-table match, empty (null) values are substituted for the right-table columns.

Exercise: There are also right outer joins and full outer joins. Try to find out what those do.

We can also join a table against itself. This is called a *self join*. As an example, suppose we wish to find all the weather records that are in the temperature range of other weather records. So we need to compare the `temp_lo` and `temp_hi` columns of each `weather` row to the `temp_lo` and `temp_hi` columns of all other `weather` rows. We can do this with the following query:

```
SELECT w1.city, w1.temp_lo AS low, w1.temp_hi AS high,
       w2.city, w2.temp_lo AS low, w2.temp_hi AS high
FROM weather w1 JOIN weather w2
ON w1.temp_lo < w2.temp_lo AND w1.temp_hi > w2.temp_hi;
```

city	low	high	city	low	high
San Francisco	43	57	San Francisco	46	50
Hayward	37	54	San Francisco	46	50

(2 rows)

Here we have relabeled the `weather` table as `w1` and `w2` to be able to distinguish the left and right side of the join. You can also use these kinds of aliases in other queries to save some typing, e.g.:

```
SELECT *
FROM weather w JOIN cities c ON w.city = c.name;
```

You will encounter this style of abbreviating quite frequently.

2.7. Aggregate Functions

Like most other relational database products, Postgres Pro supports *aggregate functions*. An aggregate function computes a single result from multiple input rows. For example, there are aggregates to compute the `count`, `sum`, `avg` (average), `max` (maximum) and `min` (minimum) over a set of rows.

As an example, we can find the highest low-temperature reading anywhere with:

```
SELECT max(temp_lo) FROM weather;
```

max
46

(1 row)

If we wanted to know what city (or cities) that reading occurred in, we might try:

```
SELECT city FROM weather WHERE temp_lo = max(temp_lo);
```

WRONG

but this will not work since the aggregate `max` cannot be used in the `WHERE` clause. (This restriction exists because the `WHERE` clause determines which rows will be included in the aggregate calculation; so obviously it has to be evaluated before aggregate functions are computed.) However, as is often the case the query can be restated to accomplish the desired result, here by using a *subquery*:

```
SELECT city FROM weather
WHERE temp_lo = (SELECT max(temp_lo) FROM weather);
```

```
city
-----
San Francisco
(1 row)
```

This is OK because the subquery is an independent computation that computes its own aggregate separately from what is happening in the outer query.

Aggregates are also very useful in combination with `GROUP BY` clauses. For example, we can get the number of readings and the maximum low temperature observed in each city with:

```
SELECT city, count(*), max(temp_lo)
FROM weather
GROUP BY city;
```

```
city      | count | max
-----+-----+-----
Hayward   |      1 |   37
San Francisco |      2 |   46
(2 rows)
```

which gives us one output row per city. Each aggregate result is computed over the table rows matching that city. We can filter these grouped rows using `HAVING`:

```
SELECT city, count(*), max(temp_lo)
FROM weather
GROUP BY city
HAVING max(temp_lo) < 40;
```

```
city      | count | max
-----+-----+-----
Hayward   |      1 |   37
(1 row)
```

which gives us the same results for only the cities that have all `temp_lo` values below 40. Finally, if we only care about cities whose names begin with “S”, we might do:

```
SELECT city, count(*), max(temp_lo)
FROM weather
WHERE city LIKE 'S%'
```

1

```
GROUP BY city;
```

```
city      | count | max
-----+-----+-----
San Francisco |      2 |   46
(1 row)
```

1 The `LIKE` operator does pattern matching and is explained in [Section 9.7](#).

It is important to understand the interaction between aggregates and SQL's `WHERE` and `HAVING` clauses. The fundamental difference between `WHERE` and `HAVING` is this: `WHERE` selects input rows before groups

and aggregates are computed (thus, it controls which rows go into the aggregate computation), whereas `HAVING` selects group rows after groups and aggregates are computed. Thus, the `WHERE` clause must not contain aggregate functions; it makes no sense to try to use an aggregate to determine which rows will be inputs to the aggregates. On the other hand, the `HAVING` clause always contains aggregate functions. (Strictly speaking, you are allowed to write a `HAVING` clause that doesn't use aggregates, but it's seldom useful. The same condition could be used more efficiently at the `WHERE` stage.)

In the previous example, we can apply the city name restriction in `WHERE`, since it needs no aggregate. This is more efficient than adding the restriction to `HAVING`, because we avoid doing the grouping and aggregate calculations for all rows that fail the `WHERE` check.

Another way to select the rows that go into an aggregate computation is to use `FILTER`, which is a per-aggregate option:

```
SELECT city, count(*) FILTER (WHERE temp_lo < 45), max(temp_lo)
FROM weather
GROUP BY city;
```

city	count	max
Hayward	1	37
San Francisco	1	46

(2 rows)

`FILTER` is much like `WHERE`, except that it removes rows only from the input of the particular aggregate function that it is attached to. Here, the `count` aggregate counts only rows with `temp_lo` below 45; but the `max` aggregate is still applied to all rows, so it still finds the reading of 46.

2.8. Updates

You can update existing rows using the `UPDATE` command. Suppose you discover the temperature readings are all off by 2 degrees after November 28. You can correct the data as follows:

```
UPDATE weather
SET temp_hi = temp_hi - 2, temp_lo = temp_lo - 2
WHERE date > '1994-11-28';
```

Look at the new state of the data:

```
SELECT * FROM weather;
```

city	temp_lo	temp_hi	prcp	date
San Francisco	46	50	0.25	1994-11-27
San Francisco	41	55	0	1994-11-29
Hayward	35	52		1994-11-29

(3 rows)

2.9. Deletions

Rows can be removed from a table using the `DELETE` command. Suppose you are no longer interested in the weather of Hayward. Then you can do the following to delete those rows from the table:

```
DELETE FROM weather WHERE city = 'Hayward';
```

All weather records belonging to Hayward are removed.

```
SELECT * FROM weather;
```

city	temp_lo	temp_hi	prcp	date
San Francisco	46	50	0.25	1994-11-27
San Francisco	41	55	0	1994-11-29

(2 rows)

One should be wary of statements of the form

```
DELETE FROM tablename;
```

Without a qualification, `DELETE` will remove *all* rows from the given table, leaving it empty. The system will not request confirmation before doing this!

Chapter 3. Advanced Features

3.1. Introduction

In the previous chapter we have covered the basics of using SQL to store and access your data in Postgres Pro. We will now discuss some more advanced features of SQL that simplify management and prevent loss or corruption of your data. Finally, we will look at some Postgres Pro extensions.

This chapter will on occasion refer to examples found in [Chapter 2](#) to change or improve them, so it will be useful to have read that chapter. Some examples from this chapter can also be found in `advanced.sql` in the tutorial directory. This file also contains some sample data to load, which is not repeated here. (Refer to [Section 2.1](#) for how to use the file.)

3.2. Views

Refer back to the queries in [Section 2.6](#). Suppose the combined listing of weather records and city location is of particular interest to your application, but you do not want to type the query each time you need it. You can create a *view* over the query, which gives a name to the query that you can refer to like an ordinary table:

```
CREATE VIEW myview AS
    SELECT name, temp_lo, temp_hi, prcp, date, location
       FROM weather, cities
      WHERE city = name;

SELECT * FROM myview;
```

Making liberal use of views is a key aspect of good SQL database design. Views allow you to encapsulate the details of the structure of your tables, which might change as your application evolves, behind consistent interfaces.

Views can be used in almost any place a real table can be used. Building views upon other views is not uncommon.

3.3. Foreign Keys

Recall the `weather` and `cities` tables from [Chapter 2](#). Consider the following problem: You want to make sure that no one can insert rows in the `weather` table that do not have a matching entry in the `cities` table. This is called maintaining the *referential integrity* of your data. In simplistic database systems this would be implemented (if at all) by first looking at the `cities` table to check if a matching record exists, and then inserting or rejecting the new `weather` records. This approach has a number of problems and is very inconvenient, so Postgres Pro can do this for you.

The new declaration of the tables would look like this:

```
CREATE TABLE cities (
    name      varchar(80) primary key,
    location  point
);

CREATE TABLE weather (
    city      varchar(80) references cities(name),
    temp_lo   int,
    temp_hi   int,
    prcp      real,
    date      date
);
```

Now try inserting an invalid record:

```
INSERT INTO weather VALUES ('Berkeley', 45, 53, 0.0, '1994-11-28');
ERROR:  insert or update on table "weather" violates foreign key constraint
"weather_city_fkey"
DETAIL:  Key (city)=(Berkeley) is not present in table "cities".
```

The behavior of foreign keys can be finely tuned to your application. We will not go beyond this simple example in this tutorial, but just refer you to [Chapter 5](#) for more information. Making correct use of foreign keys will definitely improve the quality of your database applications, so you are strongly encouraged to learn about them.

3.4. Transactions

Transactions are a fundamental concept of all database systems. The essential point of a transaction is that it bundles multiple steps into a single, all-or-nothing operation. The intermediate states between the steps are not visible to other concurrent transactions, and if some failure occurs that prevents the transaction from completing, then none of the steps affect the database at all.

For example, consider a bank database that contains balances for various customer accounts, as well as total deposit balances for branches. Suppose that we want to record a payment of \$100.00 from Alice's account to Bob's account. Simplifying outrageously, the SQL commands for this might look like:

```
UPDATE accounts SET balance = balance - 100.00
  WHERE name = 'Alice';
UPDATE branches SET balance = balance - 100.00
  WHERE name = (SELECT branch_name FROM accounts WHERE name = 'Alice');
UPDATE accounts SET balance = balance + 100.00
  WHERE name = 'Bob';
UPDATE branches SET balance = balance + 100.00
  WHERE name = (SELECT branch_name FROM accounts WHERE name = 'Bob');
```

The details of these commands are not important here; the important point is that there are several separate updates involved to accomplish this rather simple operation. Our bank's officers will want to be assured that either all these updates happen, or none of them happen. It would certainly not do for a system failure to result in Bob receiving \$100.00 that was not debited from Alice. Nor would Alice long remain a happy customer if she was debited without Bob being credited. We need a guarantee that if something goes wrong partway through the operation, none of the steps executed so far will take effect. Grouping the updates into a *transaction* gives us this guarantee. A transaction is said to be *atomic*: from the point of view of other transactions, it either happens completely or not at all.

We also want a guarantee that once a transaction is completed and acknowledged by the database system, it has indeed been permanently recorded and won't be lost even if a crash ensues shortly thereafter. For example, if we are recording a cash withdrawal by Bob, we do not want any chance that the debit to his account will disappear in a crash just after he walks out the bank door. A transactional database guarantees that all the updates made by a transaction are logged in permanent storage (i.e., on disk) before the transaction is reported complete.

Another important property of transactional databases is closely related to the notion of atomic updates: when multiple transactions are running concurrently, each one should not be able to see the incomplete changes made by others. For example, if one transaction is busy totalling all the branch balances, it would not do for it to include the debit from Alice's branch but not the credit to Bob's branch, nor vice versa. So transactions must be all-or-nothing not only in terms of their permanent effect on the database, but also in terms of their visibility as they happen. The updates made so far by an open transaction are invisible to other transactions until the transaction completes, whereupon all the updates become visible simultaneously.

In Postgres Pro, a transaction is set up by surrounding the SQL commands of the transaction with `BEGIN` and `COMMIT` commands. So our banking transaction would actually look like:

```
BEGIN;
UPDATE accounts SET balance = balance - 100.00
```

```
WHERE name = 'Alice';
-- etc etc
COMMIT;
```

If, partway through the transaction, we decide we do not want to commit (perhaps we just noticed that Alice's balance went negative), we can issue the command `ROLLBACK` instead of `COMMIT`, and all our updates so far will be canceled.

Postgres Pro actually treats every SQL statement as being executed within a transaction. If you do not issue a `BEGIN` command, then each individual statement has an implicit `BEGIN` and (if successful) `COMMIT` wrapped around it. A group of statements surrounded by `BEGIN` and `COMMIT` is sometimes called a *transaction block*.

Note

Some client libraries issue `BEGIN` and `COMMIT` commands automatically, so that you might get the effect of transaction blocks without asking. Check the documentation for the interface you are using.

It's possible to control the statements in a transaction in a more granular fashion through the use of *savepoints*. Savepoints allow you to selectively discard parts of the transaction, while committing the rest. After defining a savepoint with `SAVEPOINT`, you can if needed roll back to the savepoint with `ROLLBACK TO`. All the transaction's database changes between defining the savepoint and rolling back to it are discarded, but changes earlier than the savepoint are kept.

After rolling back to a savepoint, it continues to be defined, so you can roll back to it several times. Conversely, if you are sure you won't need to roll back to a particular savepoint again, it can be released, so the system can free some resources. Keep in mind that either releasing or rolling back to a savepoint will automatically release all savepoints that were defined after it.

All this is happening within the transaction block, so none of it is visible to other database sessions. When and if you commit the transaction block, the committed actions become visible as a unit to other sessions, while the rolled-back actions never become visible at all.

Remembering the bank database, suppose we debit \$100.00 from Alice's account, and credit Bob's account, only to find later that we should have credited Wally's account. We could do it using savepoints like this:

```
BEGIN;
UPDATE accounts SET balance = balance - 100.00
  WHERE name = 'Alice';
SAVEPOINT my_savepoint;
UPDATE accounts SET balance = balance + 100.00
  WHERE name = 'Bob';
-- oops ... forget that and use Wally's account
ROLLBACK TO my_savepoint;
UPDATE accounts SET balance = balance + 100.00
  WHERE name = 'Wally';
COMMIT;
```

This example is, of course, oversimplified, but there's a lot of control possible in a transaction block through the use of savepoints. Moreover, `ROLLBACK TO` is the only way to regain control of a transaction block that was put in aborted state by the system due to an error, short of rolling it back completely and starting again.

3.5. Window Functions

A *window function* performs a calculation across a set of table rows that are somehow related to the current row. This is comparable to the type of calculation that can be done with an aggregate function.

However, window functions do not cause rows to become grouped into a single output row like non-window aggregate calls would. Instead, the rows retain their separate identities. Behind the scenes, the window function is able to access more than just the current row of the query result.

Here is an example that shows how to compare each employee's salary with the average salary in his or her department:

```
SELECT depname, empno, salary, avg(salary) OVER (PARTITION BY depname) FROM empsalary;
```

depname	empno	salary	avg
develop	11	5200	5020.0000000000000000
develop	7	4200	5020.0000000000000000
develop	9	4500	5020.0000000000000000
develop	8	6000	5020.0000000000000000
develop	10	5200	5020.0000000000000000
personnel	5	3500	3700.0000000000000000
personnel	2	3900	3700.0000000000000000
sales	3	4800	4866.6666666666666667
sales	1	5000	4866.6666666666666667
sales	4	4800	4866.6666666666666667

(10 rows)

The first three output columns come directly from the table `empsalary`, and there is one output row for each row in the table. The fourth column represents an average taken across all the table rows that have the same `depname` value as the current row. (This actually is the same function as the non-window `avg` aggregate, but the `OVER` clause causes it to be treated as a window function and computed across the window frame.)

A window function call always contains an `OVER` clause directly following the window function's name and argument(s). This is what syntactically distinguishes it from a normal function or non-window aggregate. The `OVER` clause determines exactly how the rows of the query are split up for processing by the window function. The `PARTITION BY` clause within `OVER` divides the rows into groups, or partitions, that share the same values of the `PARTITION BY` expression(s). For each row, the window function is computed across the rows that fall into the same partition as the current row.

You can also control the order in which rows are processed by window functions using `ORDER BY` within `OVER`. (The window `ORDER BY` does not even have to match the order in which the rows are output.) Here is an example:

```
SELECT depname, empno, salary,
       rank() OVER (PARTITION BY depname ORDER BY salary DESC)
FROM empsalary;
```

depname	empno	salary	rank
develop	8	6000	1
develop	10	5200	2
develop	11	5200	2
develop	9	4500	4
develop	7	4200	5
personnel	2	3900	1
personnel	5	3500	2
sales	1	5000	1
sales	4	4800	2
sales	3	4800	2

(10 rows)

As shown here, the `rank` function produces a numerical rank for each distinct `ORDER BY` value in the current row's partition, using the order defined by the `ORDER BY` clause. `rank` needs no explicit parameter, because its behavior is entirely determined by the `OVER` clause.

The rows considered by a window function are those of the “virtual table” produced by the query's `FROM` clause as filtered by its `WHERE`, `GROUP BY`, and `HAVING` clauses if any. For example, a row removed because it does not meet the `WHERE` condition is not seen by any window function. A query can contain multiple window functions that slice up the data in different ways using different `OVER` clauses, but they all act on the same collection of rows defined by this virtual table.

We already saw that `ORDER BY` can be omitted if the ordering of rows is not important. It is also possible to omit `PARTITION BY`, in which case there is a single partition containing all rows.

There is another important concept associated with window functions: for each row, there is a set of rows within its partition called its *window frame*. Some window functions act only on the rows of the window frame, rather than of the whole partition. By default, if `ORDER BY` is supplied then the frame consists of all rows from the start of the partition up through the current row, plus any following rows that are equal to the current row according to the `ORDER BY` clause. When `ORDER BY` is omitted the default frame consists of all rows in the partition.¹ Here is an example using `sum`:

```
SELECT salary, sum(salary) OVER () FROM empsalary;
```

salary	sum
5200	47100
5000	47100
3500	47100
4800	47100
3900	47100
4200	47100
4500	47100
4800	47100
6000	47100
5200	47100

(10 rows)

Above, since there is no `ORDER BY` in the `OVER` clause, the window frame is the same as the partition, which for lack of `PARTITION BY` is the whole table; in other words each `sum` is taken over the whole table and so we get the same result for each output row. But if we add an `ORDER BY` clause, we get very different results:

```
SELECT salary, sum(salary) OVER (ORDER BY salary) FROM empsalary;
```

salary	sum
3500	3500
3900	7400
4200	11600
4500	16100
4800	25700
4800	25700
5000	30700
5200	41100
5200	41100
6000	47100

(10 rows)

Here the `sum` is taken from the first (lowest) salary up through the current one, including any duplicates of the current one (notice the results for the duplicated salaries).

Window functions are permitted only in the `SELECT` list and the `ORDER BY` clause of the query. They are forbidden elsewhere, such as in `GROUP BY`, `HAVING` and `WHERE` clauses. This is because they logically execute after the processing of those clauses. Also, window functions execute after non-window aggregate

¹ There are options to define the window frame in other ways, but this tutorial does not cover them. See [Section 4.2.8](#) for details.

functions. This means it is valid to include an aggregate function call in the arguments of a window function, but not vice versa.

If there is a need to filter or group rows after the window calculations are performed, you can use a sub-select. For example:

```
SELECT depname, empno, salary, enroll_date
FROM
  (SELECT depname, empno, salary, enroll_date,
         rank() OVER (PARTITION BY depname ORDER BY salary DESC, empno) AS pos
   FROM empsalary
  ) AS ss
WHERE pos < 3;
```

The above query only shows the rows from the inner query having `rank` less than 3.

When a query involves multiple window functions, it is possible to write out each one with a separate `OVER` clause, but this is duplicative and error-prone if the same windowing behavior is wanted for several functions. Instead, each windowing behavior can be named in a `WINDOW` clause and then referenced in `OVER`. For example:

```
SELECT sum(salary) OVER w, avg(salary) OVER w
FROM empsalary
WINDOW w AS (PARTITION BY depname ORDER BY salary DESC);
```

More details about window functions can be found in [Section 4.2.8](#), [Section 9.22](#), [Section 7.2.5](#), and the [SELECT](#) reference page.

3.6. Inheritance

Inheritance is a concept from object-oriented databases. It opens up interesting new possibilities of database design. It has some limitations in the distributed systems, described for each command (see [CREATE TABLE](#) and [ALTER TABLE](#)).

Let's create two tables: A table `cities` and a table `capitals`. Naturally, capitals are also cities, so you want some way to show the capitals implicitly when you list all cities. If you're really clever you might invent some scheme like this:

```
CREATE TABLE capitals (
  name      text,
  population real,
  elevation  int,    -- (in ft)
  state      char(2)
);

CREATE TABLE non_capitals (
  name      text,
  population real,
  elevation  int    -- (in ft)
);

CREATE VIEW cities AS
  SELECT name, population, elevation FROM capitals
  UNION
  SELECT name, population, elevation FROM non_capitals;
```

This works OK as far as querying goes, but it gets ugly when you need to update several rows, for one thing.

A better solution is this:

```
CREATE TABLE cities (
```

```
name      text,
population real,
elevation int    -- (in ft)
);

CREATE TABLE capitals (
    state    char(2) UNIQUE NOT NULL
) INHERITS (cities);
```

In this case, a row of `capitals` *inherits* all columns (`name`, `population`, and `elevation`) from its *parent*, `cities`. The type of the column `name` is `text`, a native Postgres Pro type for variable length character strings. The `capitals` table has an additional column, `state`, which shows its state abbreviation. In Postgres Pro, a table can inherit from zero or more other tables.

For example, the following query finds the names of all cities, including state capitals, that are located at an elevation over 500 feet:

```
SELECT name, elevation
FROM cities
WHERE elevation > 500;
```

which returns:

name	elevation
Las Vegas	2174
Mariposa	1953
Madison	845

(3 rows)

On the other hand, the following query finds all the cities that are not state capitals and are situated at an elevation over 500 feet:

```
SELECT name, elevation
FROM ONLY cities
WHERE elevation > 500;
```

name	elevation
Las Vegas	2174
Mariposa	1953

(2 rows)

Here the `ONLY` before `cities` indicates that the query should be run over only the `cities` table, and not tables below `cities` in the inheritance hierarchy. Many of the commands that we have already discussed — `SELECT`, `UPDATE`, and `DELETE` — support this `ONLY` notation.

Note

Although inheritance is frequently useful, it has not been integrated with unique constraints or foreign keys, which limits its usefulness. See [Section 5.11](#) for more detail.

3.7. Conclusion

Postgres Pro has many features not touched upon in this tutorial introduction, which has been oriented toward newer users of SQL. These features are discussed in more detail in the remainder of this book.

If you feel you need more introductory material, please visit the PostgreSQL [web site](#) for links to more resources.

Part II. The SQL Language

This part describes the use of the SQL language in Postgres Pro. We start with describing the general syntax of SQL, then how to create tables, how to populate the database, and how to query it. The middle part lists the available data types and functions for use in SQL commands. Lastly, we address several aspects of importance for tuning a database.

The information is arranged so that a novice user can follow it from start to end and gain a full understanding of the topics without having to refer forward too many times. The chapters are intended to be self-contained, so that advanced users can read the chapters individually as they choose. The information is presented in narrative form with topical units. Readers looking for a complete description of a particular command are encouraged to review the [Part VI](#).

Readers should know how to connect to a Postgres Pro database and issue SQL commands. Readers that are unfamiliar with these issues are encouraged to read [Part I](#) first. SQL commands are typically entered using the Postgres Pro interactive terminal `psql`, but other programs that have similar functionality can be used as well.

Chapter 4. SQL Syntax

This chapter describes the syntax of SQL. It forms the foundation for understanding the following chapters which will go into detail about how SQL commands are applied to define and modify data.

We also advise users who are already familiar with SQL to read this chapter carefully because it contains several rules and concepts that are implemented inconsistently among SQL databases or that are specific to Postgres Pro.

4.1. Lexical Structure

SQL input consists of a sequence of *commands*. A command is composed of a sequence of *tokens*, terminated by a semicolon (“;”). The end of the input stream also terminates a command. Which tokens are valid depends on the syntax of the particular command.

A token can be a *key word*, an *identifier*, a *quoted identifier*, a *literal* (or constant), or a special character symbol. Tokens are normally separated by whitespace (space, tab, newline), but need not be if there is no ambiguity (which is generally only the case if a special character is adjacent to some other token type).

For example, the following is (syntactically) valid SQL input:

```
SELECT * FROM MY_TABLE;
UPDATE MY_TABLE SET A = 5;
INSERT INTO MY_TABLE VALUES (3, 'hi there');
```

This is a sequence of three commands, one per line (although this is not required; more than one command can be on a line, and commands can usefully be split across lines).

Additionally, *comments* can occur in SQL input. They are not tokens, they are effectively equivalent to whitespace.

The SQL syntax is not very consistent regarding what tokens identify commands and which are operands or parameters. The first few tokens are generally the command name, so in the above example we would usually speak of a “SELECT”, an “UPDATE”, and an “INSERT” command. But for instance the `UPDATE` command always requires a `SET` token to appear in a certain position, and this particular variation of `INSERT` also requires a `VALUES` in order to be complete. The precise syntax rules for each command are described in [Part VI](#).

4.1.1. Identifiers and Key Words

Tokens such as `SELECT`, `UPDATE`, or `VALUES` in the example above are examples of *key words*, that is, words that have a fixed meaning in the SQL language. The tokens `MY_TABLE` and `A` are examples of *identifiers*. They identify names of tables, columns, or other database objects, depending on the command they are used in. Therefore they are sometimes simply called “names”. Key words and identifiers have the same lexical structure, meaning that one cannot know whether a token is an identifier or a key word without knowing the language. A complete list of key words can be found in [Appendix C](#).

SQL identifiers and key words must begin with a letter (a-z, but also letters with diacritical marks and non-Latin letters) or an underscore (_). Subsequent characters in an identifier or key word can be letters, underscores, digits (0-9), or dollar signs (\$). Note that dollar signs are not allowed in identifiers according to the letter of the SQL standard, so their use might render applications less portable. The SQL standard will not define a key word that contains digits or starts or ends with an underscore, so identifiers of this form are safe against possible conflict with future extensions of the standard.

The system uses no more than `NAMEDATALEN-1` bytes of an identifier; longer names can be written in commands, but they will be truncated. By default, `NAMEDATALEN` is 64 so the maximum identifier length is 63 bytes.

Key words and unquoted identifiers are case-insensitive. Therefore:

```
UPDATE MY_TABLE SET A = 5;
```

can equivalently be written as:

```
uPDaTE my_Table SeT a = 5;
```

A convention often used is to write key words in upper case and names in lower case, e.g.:

```
UPDATE my_table SET a = 5;
```

There is a second kind of identifier: the *delimited identifier* or *quoted identifier*. It is formed by enclosing an arbitrary sequence of characters in double-quotes (`"`). A delimited identifier is always an identifier, never a key word. So `"select"` could be used to refer to a column or table named `select`, whereas an unquoted `select` would be taken as a key word and would therefore provoke a parse error when used where a table or column name is expected. The example can be written with quoted identifiers like this:

```
UPDATE "my_table" SET "a" = 5;
```

Quoted identifiers can contain any character, except the character with code zero. (To include a double quote, write two double quotes.) This allows constructing table or column names that would otherwise not be possible, such as ones containing spaces or ampersands. The length limitation still applies.

Quoting an identifier also makes it case-sensitive, whereas unquoted names are always folded to lower case. For example, the identifiers `FOO`, `foo`, and `"foo"` are considered the same by Postgres Pro, but `"Foo"` and `"FOO"` are different from these three and each other. (The folding of unquoted names to lower case in Postgres Pro is incompatible with the SQL standard, which says that unquoted names should be folded to upper case. Thus, `foo` should be equivalent to `"FOO"` not `"foo"` according to the standard. If you want to write portable applications you are advised to always quote a particular name or never quote it.)

A variant of quoted identifiers allows including escaped Unicode characters identified by their code points. This variant starts with `U&` (upper or lower case `U` followed by ampersand) immediately before the opening double quote, without any spaces in between, for example `U&"foo"`. (Note that this creates an ambiguity with the operator `&`. Use spaces around the operator to avoid this problem.) Inside the quotes, Unicode characters can be specified in escaped form by writing a backslash followed by the four-digit hexadecimal code point number or alternatively a backslash followed by a plus sign followed by a six-digit hexadecimal code point number. For example, the identifier `"data"` could be written as

```
U&"d\0061t\+000061"
```

The following less trivial example writes the Russian word `"slon"` (elephant) in Cyrillic letters:

```
U&"\0441\043B\043E\043D"
```

If a different escape character than backslash is desired, it can be specified using the `UESCAPE` clause after the string, for example:

```
U&"d!0061t!+000061" UESCAPE '!'
```

The escape character can be any single character other than a hexadecimal digit, the plus sign, a single quote, a double quote, or a whitespace character. Note that the escape character is written in single quotes, not double quotes, after `UESCAPE`.

To include the escape character in the identifier literally, write it twice.

Either the 4-digit or the 6-digit escape form can be used to specify UTF-16 surrogate pairs to compose characters with code points larger than `U+FFFF`, although the availability of the 6-digit form technically makes this unnecessary. (Surrogate pairs are not stored directly, but are combined into a single code point.)

If the server encoding is not UTF-8, the Unicode code point identified by one of these escape sequences is converted to the actual server encoding; an error is reported if that's not possible.

4.1.2. Constants

There are three kinds of *implicitly-typed constants* in Postgres Pro: strings, bit strings, and numbers. Constants can also be specified with explicit types, which can enable more accurate representation and more efficient handling by the system. These alternatives are discussed in the following subsections.

4.1.2.1. String Constants

A string constant in SQL is an arbitrary sequence of characters bounded by single quotes ('), for example 'This is a string'. To include a single-quote character within a string constant, write two adjacent single quotes, e.g., 'Dianne''s horse'. Note that this is *not* the same as a double-quote character (").

Two string constants that are only separated by whitespace *with at least one newline* are concatenated and effectively treated as if the string had been written as one constant. For example:

```
SELECT 'foo'
'bar';
```

is equivalent to:

```
SELECT 'foobar';
```

but:

```
SELECT 'foo'      'bar';
```

is not valid syntax. (This slightly bizarre behavior is specified by SQL; Postgres Pro is following the standard.)

4.1.2.2. String Constants with C-Style Escapes

Postgres Pro also accepts “escape” string constants, which are an extension to the SQL standard. An escape string constant is specified by writing the letter E (upper or lower case) just before the opening single quote, e.g., E'foo'. (When continuing an escape string constant across lines, write E only before the first opening quote.) Within an escape string, a backslash character (\) begins a C-like *backslash escape* sequence, in which the combination of backslash and following character(s) represent a special byte value, as shown in [Table 4.1](#).

Table 4.1. Backslash Escape Sequences

Backslash Escape Sequence	Interpretation
\b	backspace
\f	form feed
\n	newline
\r	carriage return
\t	tab
\o, \oo, \ooo (o = 0-7)	octal byte value
\xh, \xhh (h = 0-9, A-F)	hexadecimal byte value
\uxxxx, \Uxxxxxxxx (x = 0-9, A-F)	16 or 32-bit hexadecimal Unicode character value

Any other character following a backslash is taken literally. Thus, to include a backslash character, write two backslashes (\\). Also, a single quote can be included in an escape string by writing \', in addition to the normal way of ' '.

It is your responsibility that the byte sequences you create, especially when using the octal or hexadecimal escapes, compose valid characters in the server character set encoding. A useful alternative is to use Unicode escapes or the alternative Unicode escape syntax, explained in [Section 4.1.2.3](#); then the server will check that the character conversion is possible.

Caution

If the configuration parameter [standard_conforming_strings](#) is `off`, then Postgres Pro recognizes backslash escapes in both regular and escape string constants. However, as of PostgreSQL 9.1, the default is `on`, meaning that backslash escapes are recognized only in escape string constants. This

behavior is more standards-compliant, but might break applications which rely on the historical behavior, where backslash escapes were always recognized. As a workaround, you can set this parameter to `off`, but it is better to migrate away from using backslash escapes. If you need to use a backslash escape to represent a special character, write the string constant with an `E`.

In addition to `standard_conforming_strings`, the configuration parameters `escape_string_warning` and `backslash_quote` govern treatment of backslashes in string constants.

The character with the code zero cannot be in a string constant.

4.1.2.3. String Constants with Unicode Escapes

Postgres Pro also supports another type of escape syntax for strings that allows specifying arbitrary Unicode characters by code point. A Unicode escape string constant starts with `U&` (upper or lower case letter `U` followed by ampersand) immediately before the opening quote, without any spaces in between, for example `U&'foo'`. (Note that this creates an ambiguity with the operator `&`. Use spaces around the operator to avoid this problem.) Inside the quotes, Unicode characters can be specified in escaped form by writing a backslash followed by the four-digit hexadecimal code point number or alternatively a backslash followed by a plus sign followed by a six-digit hexadecimal code point number. For example, the string `'data'` could be written as

```
U&'d\0061t\+000061'
```

The following less trivial example writes the Russian word “слон” (elephant) in Cyrillic letters:

```
U&'\'0441\'043B\'043E\'043D'
```

If a different escape character than backslash is desired, it can be specified using the `UESCAPE` clause after the string, for example:

```
U&'d!0061t!+000061' UESCAPE '!'
```

The escape character can be any single character other than a hexadecimal digit, the plus sign, a single quote, a double quote, or a whitespace character.

To include the escape character in the string literally, write it twice.

Either the 4-digit or the 6-digit escape form can be used to specify UTF-16 surrogate pairs to compose characters with code points larger than `U+FFFF`, although the availability of the 6-digit form technically makes this unnecessary. (Surrogate pairs are not stored directly, but are combined into a single code point.)

If the server encoding is not UTF-8, the Unicode code point identified by one of these escape sequences is converted to the actual server encoding; an error is reported if that's not possible.

Also, the Unicode escape syntax for string constants only works when the configuration parameter `standard_conforming_strings` is turned on. This is because otherwise this syntax could confuse clients that parse the SQL statements to the point that it could lead to SQL injections and similar security issues. If the parameter is set to `off`, this syntax will be rejected with an error message.

4.1.2.4. Dollar-Quoted String Constants

While the standard syntax for specifying string constants is usually convenient, it can be difficult to understand when the desired string contains many single quotes, since each of those must be doubled. To allow more readable queries in such situations, Postgres Pro provides another way, called “dollar quoting”, to write string constants. A dollar-quoted string constant consists of a dollar sign (`$`), an optional “tag” of zero or more characters, another dollar sign, an arbitrary sequence of characters that makes up the string content, a dollar sign, the same tag that began this dollar quote, and a dollar sign. For example, here are two different ways to specify the string “Dianne's horse” using dollar quoting:

```
$$Dianne's horse$$
```

```
$SomeTag$Dianne's horse$SomeTag$
```

Notice that inside the dollar-quoted string, single quotes can be used without needing to be escaped. Indeed, no characters inside a dollar-quoted string are ever escaped: the string content is always written literally. Backslashes are not special, and neither are dollar signs, unless they are part of a sequence matching the opening tag.

It is possible to nest dollar-quoted string constants by choosing different tags at each nesting level. This is most commonly used in writing function definitions. For example:

```
$function$
BEGIN
    RETURN ($1 ~ $q$[\t\r\n\v\\]$q$);
END;
$function$
```

Here, the sequence `q[\t\r\n\v\\]q` represents a dollar-quoted literal string `[\t\r\n\v\\]`, which will be recognized when the function body is executed by Postgres Pro. But since the sequence does not match the outer dollar quoting delimiter `$function$`, it is just some more characters within the constant so far as the outer string is concerned.

The tag, if any, of a dollar-quoted string follows the same rules as an unquoted identifier, except that it cannot contain a dollar sign. Tags are case sensitive, so `tagString contenttag` is correct, but `TAGString contenttag` is not.

A dollar-quoted string that follows a keyword or identifier must be separated from it by whitespace; otherwise the dollar quoting delimiter would be taken as part of the preceding identifier.

Dollar quoting is not part of the SQL standard, but it is often a more convenient way to write complicated string literals than the standard-compliant single quote syntax. It is particularly useful when representing string constants inside other constants, as is often needed in procedural function definitions. With single-quote syntax, each backslash in the above example would have to be written as four backslashes, which would be reduced to two backslashes in parsing the original string constant, and then to one when the inner string constant is re-parsed during function execution.

4.1.2.5. Bit-String Constants

Bit-string constants look like regular string constants with a `B` (upper or lower case) immediately before the opening quote (no intervening whitespace), e.g., `B'1001'`. The only characters allowed within bit-string constants are `0` and `1`.

Alternatively, bit-string constants can be specified in hexadecimal notation, using a leading `x` (upper or lower case), e.g., `x'1FF'`. This notation is equivalent to a bit-string constant with four binary digits for each hexadecimal digit.

Both forms of bit-string constant can be continued across lines in the same way as regular string constants. Dollar quoting cannot be used in a bit-string constant.

4.1.2.6. Numeric Constants

Numeric constants are accepted in these general forms:

```
digits
digits. [digits] [e[+-] digits]
[digits]. digits [e[+-] digits]
digitse[+-] digits
```

where *digits* is one or more decimal digits (0 through 9). At least one digit must be before or after the decimal point, if one is used. At least one digit must follow the exponent marker (*e*), if one is present. There cannot be any spaces or other characters embedded in the constant, except for underscores, which can be used for visual grouping as described below. Note that any leading plus or minus sign is not actually considered part of the constant; it is an operator applied to the constant.

These are some examples of valid numeric constants:

```
42
3.5
4.
.001
5e2
1.925e-3
```

Additionally, non-decimal integer constants are accepted in these forms:

```
0xhexdigits
0ooctdigits
0bbindigits
```

where *hexdigits* is one or more hexadecimal digits (0-9, A-F), *octdigits* is one or more octal digits (0-7), and *bindigits* is one or more binary digits (0 or 1). Hexadecimal digits and the radix prefixes can be in upper or lower case. Note that only integers can have non-decimal forms, not numbers with fractional parts.

These are some examples of valid non-decimal integer constants:

```
0b100101
0B10011001
0o273
0O755
0x42f
0XFFFF
```

For visual grouping, underscores can be inserted between digits. These have no further effect on the value of the constant. For example:

```
1_500_000_000
0b10001000_00000000
0o_1_755
0xFFFF_FFFF
1.618_034
```

Underscores are not allowed at the start or end of a numeric constant or a group of digits (that is, immediately before or after the decimal point or the exponent marker), and more than one underscore in a row is not allowed.

A numeric constant that contains neither a decimal point nor an exponent is initially presumed to be type `integer` if its value fits in type `integer` (32 bits); otherwise it is presumed to be type `bigint` if its value fits in type `bigint` (64 bits); otherwise it is taken to be type `numeric`. Constants that contain decimal points and/or exponents are always initially presumed to be type `numeric`.

The initially assigned data type of a numeric constant is just a starting point for the type resolution algorithms. In most cases the constant will be automatically coerced to the most appropriate type depending on context. When necessary, you can force a numeric value to be interpreted as a specific data type by casting it. For example, you can force a numeric value to be treated as type `real` (`float4`) by writing:

```
REAL '1.23' -- string style
1.23::REAL -- Postgres Pro (historical) style
```

These are actually just special cases of the general casting notations discussed next.

4.1.2.7. Constants of Other Types

A constant of an *arbitrary* type can be entered using any one of the following notations:

```
type 'string'
'string'::type
```

```
CAST ( 'string' AS type )
```

The string constant's text is passed to the input conversion routine for the type called *type*. The result is a constant of the indicated type. The explicit type cast can be omitted if there is no ambiguity as to the type the constant must be (for example, when it is assigned directly to a table column), in which case it is automatically coerced.

The string constant can be written using either regular SQL notation or dollar-quoting.

It is also possible to specify a type coercion using a function-like syntax:

```
typename ( 'string' )
```

but not all type names can be used in this way; see [Section 4.2.9](#) for details.

The `::`, `CAST()`, and function-call syntaxes can also be used to specify run-time type conversions of arbitrary expressions, as discussed in [Section 4.2.9](#). To avoid syntactic ambiguity, the *type* `'string'` syntax can only be used to specify the type of a simple literal constant. Another restriction on the *type* `'string'` syntax is that it does not work for array types; use `::` or `CAST()` to specify the type of an array constant.

The `CAST()` syntax conforms to SQL. The *type* `'string'` syntax is a generalization of the standard: SQL specifies this syntax only for a few data types, but Postgres Pro allows it for all types. The syntax with `::` is historical Postgres Pro usage, as is the function-call syntax.

4.1.3. Operators

An operator name is a sequence of up to `NAMEDATALEN-1` (63 by default) characters from the following list:

+ - * / < > = ~ ! @ # % ^ & | ` ?

There are a few restrictions on operator names, however:

- `--` and `/*` cannot appear anywhere in an operator name, since they will be taken as the start of a comment.
- A multiple-character operator name cannot end in `+` or `-`, unless the name also contains at least one of these characters:

~ ! @ # % ^ & | ` ?

For example, `@-` is an allowed operator name, but `*-` is not. This restriction allows Postgres Pro to parse SQL-compliant queries without requiring spaces between tokens.

When working with non-SQL-standard operator names, you will usually need to separate adjacent operators with spaces to avoid ambiguity. For example, if you have defined a prefix operator named `@`, you cannot write `X*@Y`; you must write `X* @Y` to ensure that Postgres Pro reads it as two operator names not one.

4.1.4. Special Characters

Some characters that are not alphanumeric have a special meaning that is different from being an operator. Details on the usage can be found at the location where the respective syntax element is described. This section only exists to advise the existence and summarize the purposes of these characters.

- A dollar sign (`$`) followed by digits is used to represent a positional parameter in the body of a function definition or a prepared statement. In other contexts the dollar sign can be part of an identifier or a dollar-quoted string constant.
- Parentheses (`()`) have their usual meaning to group expressions and enforce precedence. In some cases parentheses are required as part of the fixed syntax of a particular SQL command.
- Brackets (`[]`) are used to select the elements of an array. See [Section 8.15](#) for more information on arrays.

- Commas (,) are used in some syntactical constructs to separate the elements of a list.
- The semicolon (;) terminates an SQL command. It cannot appear anywhere within a command, except within a string constant or quoted identifier.
- The colon (:) is used to select “slices” from arrays. (See [Section 8.15.](#)) In certain SQL dialects (such as Embedded SQL), the colon is used to prefix variable names.
- The asterisk (*) is used in some contexts to denote all the fields of a table row or composite value. It also has a special meaning when used as the argument of an aggregate function, namely that the aggregate does not require any explicit parameter.
- The period (.) is used in numeric constants, and to separate schema, table, and column names.

4.1.5. Comments

A comment is a sequence of characters beginning with double dashes and extending to the end of the line, e.g.:

```
-- This is a standard SQL comment
```

Alternatively, C-style block comments can be used:

```
/* multiline comment
 * with nesting: /* nested block comment */
 */
```

where the comment begins with `/*` and extends to the matching occurrence of `*/`. These block comments nest, as specified in the SQL standard but unlike C, so that one can comment out larger blocks of code that might contain existing block comments.

A comment is removed from the input stream before further syntax analysis and is effectively replaced by whitespace.

4.1.6. Operator Precedence

[Table 4.2](#) shows the precedence and associativity of the operators in Postgres Pro. Most operators have the same precedence and are left-associative. The precedence and associativity of the operators is hard-wired into the parser. Add parentheses if you want an expression with multiple operators to be parsed in some other way than what the precedence rules imply.

Table 4.2. Operator Precedence (highest to lowest)

Operator/Element	Associativity	Description
.	left	table/column name separator
::	left	Postgres Pro-style typecast
[]	left	array element selection
+ -	right	unary plus, unary minus
COLLATE	left	collation selection
AT	left	AT TIME ZONE, AT LOCAL
^	left	exponentiation
* / %	left	multiplication, division, modulo
+ -	left	addition, subtraction
(any other operator)	left	all other native and user-defined operators
BETWEEN IN LIKE ILIKE SIMILAR		range containment, set membership, string matching
< > = <= >= <>		comparison operators

Operator/Element	Associativity	Description
IS ISNULL NOTNULL		IS TRUE, IS FALSE, IS NULL, IS DISTINCT FROM, etc.
NOT	right	logical negation
AND	left	logical conjunction
OR	left	logical disjunction

Note that the operator precedence rules also apply to user-defined operators that have the same names as the built-in operators mentioned above. For example, if you define a “+” operator for some custom data type it will have the same precedence as the built-in “+” operator, no matter what yours does.

When a schema-qualified operator name is used in the `OPERATOR` syntax, as for example in:

```
SELECT 3 OPERATOR(pg_catalog.+) 4;
```

the `OPERATOR` construct is taken to have the default precedence shown in [Table 4.2](#) for “any other operator”. This is true no matter which specific operator appears inside `OPERATOR()`.

Note

PostgreSQL versions before 9.5 used slightly different operator precedence rules. In particular, `<=` `>=` and `<>` used to be treated as generic operators; `IS` tests used to have higher priority; and `NOT BETWEEN` and related constructs acted inconsistently, being taken in some cases as having the precedence of `NOT` rather than `BETWEEN`. These rules were changed for better compliance with the SQL standard and to reduce confusion from inconsistent treatment of logically equivalent constructs. In most cases, these changes will result in no behavioral change, or perhaps in “no such operator” failures which can be resolved by adding parentheses. However there are corner cases in which a query might change behavior without any parsing error being reported.

4.2. Value Expressions

Value expressions are used in a variety of contexts, such as in the target list of the `SELECT` command, as new column values in `INSERT` or `UPDATE`, or in search conditions in a number of commands. The result of a value expression is sometimes called a *scalar*, to distinguish it from the result of a table expression (which is a table). Value expressions are therefore also called *scalar expressions* (or even simply *expressions*). The expression syntax allows the calculation of values from primitive parts using arithmetic, logical, set, and other operations.

A value expression is one of the following:

- A constant or literal value
- A column reference
- A positional parameter reference, in the body of a function definition or prepared statement
- A subscripted expression
- A field selection expression
- An operator invocation
- A function call
- An aggregate expression
- A window function call
- A type cast

- A collation expression
- A scalar subquery
- An array constructor
- A row constructor
- Another value expression in parentheses (used to group subexpressions and override precedence)

In addition to this list, there are a number of constructs that can be classified as an expression but do not follow any general syntax rules. These generally have the semantics of a function or operator and are explained in the appropriate location in [Chapter 9](#). An example is the `IS NULL` clause.

We have already discussed constants in [Section 4.1.2](#). The following sections discuss the remaining options.

4.2.1. Column References

A column can be referenced in the form:

correlation.columnname

correlation is the name of a table (possibly qualified with a schema name), or an alias for a table defined by means of a `FROM` clause. The correlation name and separating dot can be omitted if the column name is unique across all the tables being used in the current query. (See also [Chapter 7](#).)

4.2.2. Positional Parameters

A positional parameter reference is used to indicate a value that is supplied externally to an SQL statement. Parameters are used in SQL function definitions and in prepared queries. Some client libraries also support specifying data values separately from the SQL command string, in which case parameters are used to refer to the out-of-line data values. The form of a parameter reference is:

\$number

For example, consider the definition of a function, `dept`, as:

```
CREATE FUNCTION dept(text) RETURNS dept
  AS $$ SELECT * FROM dept WHERE name = $1 $$
LANGUAGE SQL;
```

Here the `$1` references the value of the first function argument whenever the function is invoked.

4.2.3. Subscripts

If an expression yields a value of an array type, then a specific element of the array value can be extracted by writing

expression[subscript]

or multiple adjacent elements (an “array slice”) can be extracted by writing

expression[lower_subscript:upper_subscript]

(Here, the brackets `[]` are meant to appear literally.) Each *subscript* is itself an expression, which will be rounded to the nearest integer value.

In general the array *expression* must be parenthesized, but the parentheses can be omitted when the expression to be subscripted is just a column reference or positional parameter. Also, multiple subscripts can be concatenated when the original array is multidimensional. For example:

```
mytable.arraycolumn[4]
mytable.two_d_column[17][34]
$1[10:42]
```

```
(arrayfunction(a,b)) [42]
```

The parentheses in the last example are required. See [Section 8.15](#) for more about arrays.

4.2.4. Field Selection

If an expression yields a value of a composite type (row type), then a specific field of the row can be extracted by writing

```
expression.fieldname
```

In general the row *expression* must be parenthesized, but the parentheses can be omitted when the expression to be selected from is just a table reference or positional parameter. For example:

```
mytable.mycolumn  
$1.somecolumn  
(rowfunction(a,b)).col3
```

(Thus, a qualified column reference is actually just a special case of the field selection syntax.) An important special case is extracting a field from a table column that is of a composite type:

```
(compositecol).somefield  
(mytable.compositecol).somefield
```

The parentheses are required here to show that *compositecol* is a column name not a table name, or that *mytable* is a table name not a schema name in the second case.

You can ask for all fields of a composite value by writing *.**:

```
(compositecol).*
```

This notation behaves differently depending on context; see [Section 8.16.5](#) for details.

4.2.5. Operator Invocations

There are two possible syntaxes for an operator invocation:

```
expression operator expression (binary infix operator)  
operator expression (unary prefix operator)
```

where the *operator* token follows the syntax rules of [Section 4.1.3](#), or is one of the key words `AND`, `OR`, and `NOT`, or is a qualified operator name in the form:

```
OPERATOR(schema.operatorname)
```

Which particular operators exist and whether they are unary or binary depends on what operators have been defined by the system or the user. [Chapter 9](#) describes the built-in operators.

4.2.6. Function Calls

The syntax for a function call is the name of a function (possibly qualified with a schema name), followed by its argument list enclosed in parentheses:

```
function_name ([expression [, expression ... ]])
```

For example, the following computes the square root of 2:

```
sqrt(2)
```

The list of built-in functions is in [Chapter 9](#). Other functions can be added by the user.

When issuing queries in a database where some users mistrust other users, observe security precautions from [Section 10.3](#) when writing function calls.

The arguments can optionally have names attached. See [Section 4.3](#) for details.

Note

A function that takes a single argument of composite type can optionally be called using field-selection syntax, and conversely field selection can be written in functional style. That is, the notations `col(table)` and `table.col` are interchangeable. This behavior is not SQL-standard but is provided in Postgres Pro because it allows use of functions to emulate “computed fields”. For more information see [Section 8.16.5](#).

4.2.7. Aggregate Expressions

An *aggregate expression* represents the application of an aggregate function across the rows selected by a query. An aggregate function reduces multiple inputs to a single output value, such as the sum or average of the inputs. The syntax of an aggregate expression is one of the following:

```
aggregate_name (expression [ , ... ] [ order_by_clause ] ) [ FILTER
  ( WHERE filter_clause ) ]
aggregate_name (ALL expression [ , ... ] [ order_by_clause ] ) [ FILTER
  ( WHERE filter_clause ) ]
aggregate_name (DISTINCT expression [ , ... ] [ order_by_clause ] ) [ FILTER
  ( WHERE filter_clause ) ]
aggregate_name ( * ) [ FILTER ( WHERE filter_clause ) ]
aggregate_name ( [ expression [ , ... ] ] ) WITHIN GROUP ( order_by_clause ) [ FILTER
  ( WHERE filter_clause ) ]
```

where *aggregate_name* is a previously defined aggregate (possibly qualified with a schema name) and *expression* is any value expression that does not itself contain an aggregate expression or a window function call. The optional *order_by_clause* and *filter_clause* are described below.

The first form of aggregate expression invokes the aggregate once for each input row. The second form is the same as the first, since `ALL` is the default. The third form invokes the aggregate once for each distinct value of the expression (or distinct set of values, for multiple expressions) found in the input rows. The fourth form invokes the aggregate once for each input row; since no particular input value is specified, it is generally only useful for the `count(*)` aggregate function. The last form is used with *ordered-set* aggregate functions, which are described below.

Most aggregate functions ignore null inputs, so that rows in which one or more of the expression(s) yield null are discarded. This can be assumed to be true, unless otherwise specified, for all built-in aggregates.

For example, `count(*)` yields the total number of input rows; `count(f1)` yields the number of input rows in which `f1` is non-null, since `count` ignores nulls; and `count(distinct f1)` yields the number of distinct non-null values of `f1`.

Ordinarily, the input rows are fed to the aggregate function in an unspecified order. In many cases this does not matter; for example, `min` produces the same result no matter what order it receives the inputs in. However, some aggregate functions (such as `array_agg` and `string_agg`) produce results that depend on the ordering of the input rows. When using such an aggregate, the optional *order_by_clause* can be used to specify the desired ordering. The *order_by_clause* has the same syntax as for a query-level `ORDER BY` clause, as described in [Section 7.5](#), except that its expressions are always just expressions and cannot be output-column names or numbers. For example:

```
WITH vals (v) AS ( VALUES (1),(3),(4),(3),(2) )
SELECT array_agg(v ORDER BY v DESC) FROM vals;
   array_agg
-----
{4,3,3,2,1}
```

Since `jsonb` only keeps the last matching key, ordering of its keys can be significant:

```
WITH vals (k, v) AS ( VALUES ('key0','1'), ('key1','3'), ('key1','2') )
```

```
SELECT jsonb_object_agg(k, v ORDER BY v) FROM vals;
       jsonb_object_agg
-----
{"key0": "1", "key1": "3"}
```

When dealing with multiple-argument aggregate functions, note that the `ORDER BY` clause goes after all the aggregate arguments. For example, write this:

```
SELECT string_agg(a, ',' ORDER BY a) FROM table;
```

not this:

```
SELECT string_agg(a ORDER BY a, ',') FROM table; -- incorrect
```

The latter is syntactically valid, but it represents a call of a single-argument aggregate function with two `ORDER BY` keys (the second one being rather useless since it's a constant).

If `DISTINCT` is specified with an *order_by_clause*, `ORDER BY` expressions can only reference columns in the `DISTINCT` list. For example:

```
WITH vals (v) AS ( VALUES (1), (3), (4), (3), (2) )
SELECT array_agg(DISTINCT v ORDER BY v DESC) FROM vals;
       array_agg
-----
{4,3,2,1}
```

Placing `ORDER BY` within the aggregate's regular argument list, as described so far, is used when ordering the input rows for general-purpose and statistical aggregates, for which ordering is optional. There is a subclass of aggregate functions called *ordered-set aggregates* for which an *order_by_clause* is *required*, usually because the aggregate's computation is only sensible in terms of a specific ordering of its input rows. Typical examples of ordered-set aggregates include rank and percentile calculations. For an ordered-set aggregate, the *order_by_clause* is written inside `WITHIN GROUP (...)`, as shown in the final syntax alternative above. The expressions in the *order_by_clause* are evaluated once per input row just like regular aggregate arguments, sorted as per the *order_by_clause*'s requirements, and fed to the aggregate function as input arguments. (This is unlike the case for a non-`WITHIN GROUP` *order_by_clause*, which is not treated as argument(s) to the aggregate function.) The argument expressions preceding `WITHIN GROUP`, if any, are called *direct arguments* to distinguish them from the *aggregated arguments* listed in the *order_by_clause*. Unlike regular aggregate arguments, direct arguments are evaluated only once per aggregate call, not once per input row. This means that they can contain variables only if those variables are grouped by `GROUP BY`; this restriction is the same as if the direct arguments were not inside an aggregate expression at all. Direct arguments are typically used for things like percentile fractions, which only make sense as a single value per aggregation calculation. The direct argument list can be empty; in this case, write just `()` not `(*)`. (Postgres Pro will actually accept either spelling, but only the first way conforms to the SQL standard.)

An example of an ordered-set aggregate call is:

```
SELECT percentile_cont(0.5) WITHIN GROUP (ORDER BY income) FROM households;
       percentile_cont
-----
50489
```

which obtains the 50th percentile, or median, value of the `income` column from table `households`. Here, `0.5` is a direct argument; it would make no sense for the percentile fraction to be a value varying across rows.

If `FILTER` is specified, then only the input rows for which the *filter_clause* evaluates to true are fed to the aggregate function; other rows are discarded. For example:

```
SELECT
    count(*) AS unfiltered,
    count(*) FILTER (WHERE i < 5) AS filtered
FROM generate_series(1,10) AS s(i);
```

```
unfiltered | filtered
-----+-----
          10 |          4
(1 row)
```

The predefined aggregate functions are described in [Section 9.21](#). Other aggregate functions can be added by the user.

An aggregate expression can only appear in the result list or `HAVING` clause of a `SELECT` command. It is forbidden in other clauses, such as `WHERE`, because those clauses are logically evaluated before the results of aggregates are formed.

When an aggregate expression appears in a subquery (see [Section 4.2.11](#) and [Section 9.24](#)), the aggregate is normally evaluated over the rows of the subquery. But an exception occurs if the aggregate's arguments (and *filter_clause* if any) contain only outer-level variables: the aggregate then belongs to the nearest such outer level, and is evaluated over the rows of that query. The aggregate expression as a whole is then an outer reference for the subquery it appears in, and acts as a constant over any one evaluation of that subquery. The restriction about appearing only in the result list or `HAVING` clause applies with respect to the query level that the aggregate belongs to.

4.2.8. Window Function Calls

A *window function call* represents the application of an aggregate-like function over some portion of the rows selected by a query. Unlike non-window aggregate calls, this is not tied to grouping of the selected rows into a single output row — each row remains separate in the query output. However the window function has access to all the rows that would be part of the current row's group according to the grouping specification (`PARTITION BY` list) of the window function call. The syntax of a window function call is one of the following:

```
function_name ([expression [, expression ... ]]) [ FILTER ( WHERE filter_clause ) ]
OVER window_name
function_name ([expression [, expression ... ]]) [ FILTER ( WHERE filter_clause ) ]
OVER ( window_definition )
function_name ( * ) [ FILTER ( WHERE filter_clause ) ] OVER window_name
function_name ( * ) [ FILTER ( WHERE filter_clause ) ] OVER ( window_definition )
```

where *window_definition* has the syntax

```
[ existing_window_name ]
[ PARTITION BY expression [, ...] ]
[ ORDER BY expression [ ASC | DESC | USING operator ] [ NULLS { FIRST | LAST } ]
[, ...] ]
[ frame_clause ]
```

The optional *frame_clause* can be one of

```
{ RANGE | ROWS | GROUPS } frame_start [ frame_exclusion ]
{ RANGE | ROWS | GROUPS } BETWEEN frame_start AND frame_end [ frame_exclusion ]
```

where *frame_start* and *frame_end* can be one of

```
UNBOUNDED PRECEDING
offset PRECEDING
CURRENT ROW
offset FOLLOWING
UNBOUNDED FOLLOWING
```

and *frame_exclusion* can be one of

```
EXCLUDE CURRENT ROW
EXCLUDE GROUP
EXCLUDE TIES
EXCLUDE NO OTHERS
```

Here, *expression* represents any value expression that does not itself contain window function calls.

window_name is a reference to a named window specification defined in the query's `WINDOW` clause. Alternatively, a full *window_definition* can be given within parentheses, using the same syntax as for defining a named window in the `WINDOW` clause; see the [SELECT](#) reference page for details. It's worth pointing out that `OVER wname` is not exactly equivalent to `OVER (wname ...)`; the latter implies copying and modifying the window definition, and will be rejected if the referenced window specification includes a frame clause.

The `PARTITION BY` clause groups the rows of the query into *partitions*, which are processed separately by the window function. `PARTITION BY` works similarly to a query-level `GROUP BY` clause, except that its expressions are always just expressions and cannot be output-column names or numbers. Without `PARTITION BY`, all rows produced by the query are treated as a single partition. The `ORDER BY` clause determines the order in which the rows of a partition are processed by the window function. It works similarly to a query-level `ORDER BY` clause, but likewise cannot use output-column names or numbers. Without `ORDER BY`, rows are processed in an unspecified order.

The *frame_clause* specifies the set of rows constituting the *window frame*, which is a subset of the current partition, for those window functions that act on the frame instead of the whole partition. The set of rows in the frame can vary depending on which row is the current row. The frame can be specified in `RANGE`, `ROWS` or `GROUPS` mode; in each case, it runs from the *frame_start* to the *frame_end*. If *frame_end* is omitted, the end defaults to `CURRENT ROW`.

A *frame_start* of `UNBOUNDED PRECEDING` means that the frame starts with the first row of the partition, and similarly a *frame_end* of `UNBOUNDED FOLLOWING` means that the frame ends with the last row of the partition.

In `RANGE` or `GROUPS` mode, a *frame_start* of `CURRENT ROW` means the frame starts with the current row's first *peer* row (a row that the window's `ORDER BY` clause sorts as equivalent to the current row), while a *frame_end* of `CURRENT ROW` means the frame ends with the current row's last peer row. In `ROWS` mode, `CURRENT ROW` simply means the current row.

In the *offset* `PRECEDING` and *offset* `FOLLOWING` frame options, the *offset* must be an expression not containing any variables, aggregate functions, or window functions. The meaning of the *offset* depends on the frame mode:

- In `ROWS` mode, the *offset* must yield a non-null, non-negative integer, and the option means that the frame starts or ends the specified number of rows before or after the current row.
- In `GROUPS` mode, the *offset* again must yield a non-null, non-negative integer, and the option means that the frame starts or ends the specified number of *peer groups* before or after the current row's peer group, where a peer group is a set of rows that are equivalent in the `ORDER BY` ordering. (There must be an `ORDER BY` clause in the window definition to use `GROUPS` mode.)
- In `RANGE` mode, these options require that the `ORDER BY` clause specify exactly one column. The *offset* specifies the maximum difference between the value of that column in the current row and its value in preceding or following rows of the frame. The data type of the *offset* expression varies depending on the data type of the ordering column. For numeric ordering columns it is typically of the same type as the ordering column, but for datetime ordering columns it is an interval. For example, if the ordering column is of type `date` or `timestamp`, one could write `RANGE BETWEEN '1 day' PRECEDING AND '10 days' FOLLOWING`. The *offset* is still required to be non-null and non-negative, though the meaning of “non-negative” depends on its data type.

In any case, the distance to the end of the frame is limited by the distance to the end of the partition, so that for rows near the partition ends the frame might contain fewer rows than elsewhere.

Notice that in both `ROWS` and `GROUPS` mode, `0 PRECEDING` and `0 FOLLOWING` are equivalent to `CURRENT ROW`. This normally holds in `RANGE` mode as well, for an appropriate data-type-specific meaning of “zero”.

The *frame_exclusion* option allows rows around the current row to be excluded from the frame, even if they would be included according to the frame start and frame end options. `EXCLUDE CURRENT ROW`

excludes the current row from the frame. `EXCLUDE GROUP` excludes the current row and its ordering peers from the frame. `EXCLUDE TIES` excludes any peers of the current row from the frame, but not the current row itself. `EXCLUDE NO OTHERS` simply specifies explicitly the default behavior of not excluding the current row or its peers.

The default framing option is `RANGE UNBOUNDED PRECEDING`, which is the same as `RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW`. With `ORDER BY`, this sets the frame to be all rows from the partition start up through the current row's last `ORDER BY` peer. Without `ORDER BY`, this means all rows of the partition are included in the window frame, since all rows become peers of the current row.

Restrictions are that *frame_start* cannot be `UNBOUNDED FOLLOWING`, *frame_end* cannot be `UNBOUNDED PRECEDING`, and the *frame_end* choice cannot appear earlier in the above list of *frame_start* and *frame_end* options than the *frame_start* choice does — for example `RANGE BETWEEN CURRENT ROW AND offset PRECEDING` is not allowed. But, for example, `ROWS BETWEEN 7 PRECEDING AND 8 PRECEDING` is allowed, even though it would never select any rows.

If `FILTER` is specified, then only the input rows for which the *filter_clause* evaluates to true are fed to the window function; other rows are discarded. Only window functions that are aggregates accept a `FILTER` clause.

The built-in window functions are described in [Table 9.65](#). Other window functions can be added by the user. Also, any built-in or user-defined general-purpose or statistical aggregate can be used as a window function. (Ordered-set and hypothetical-set aggregates cannot presently be used as window functions.)

The syntaxes using `*` are used for calling parameter-less aggregate functions as window functions, for example `count(*) OVER (PARTITION BY x ORDER BY y)`. The asterisk (`*`) is customarily not used for window-specific functions. Window-specific functions do not allow `DISTINCT` or `ORDER BY` to be used within the function argument list.

Window function calls are permitted only in the `SELECT` list and the `ORDER BY` clause of the query.

More information about window functions can be found in [Section 3.5](#), [Section 9.22](#), and [Section 7.2.5](#).

4.2.9. Type Casts

A type cast specifies a conversion from one data type to another. Postgres Pro accepts two equivalent syntaxes for type casts:

```
CAST ( expression AS type )
expression::type
```

The `CAST` syntax conforms to SQL; the syntax with `::` is historical Postgres Pro usage.

When a cast is applied to a value expression of a known type, it represents a run-time type conversion. The cast will succeed only if a suitable type conversion operation has been defined. Notice that this is subtly different from the use of casts with constants, as shown in [Section 4.1.2.7](#). A cast applied to an unadorned string literal represents the initial assignment of a type to a literal constant value, and so it will succeed for any type (if the contents of the string literal are acceptable input syntax for the data type).

An explicit type cast can usually be omitted if there is no ambiguity as to the type that a value expression must produce (for example, when it is assigned to a table column); the system will automatically apply a type cast in such cases. However, automatic casting is only done for casts that are marked “OK to apply implicitly” in the system catalogs. Other casts must be invoked with explicit casting syntax. This restriction is intended to prevent surprising conversions from being applied silently.

It is also possible to specify a type cast using a function-like syntax:

```
typename ( expression )
```

However, this only works for types whose names are also valid as function names. For example, `double precision` cannot be used this way, but the equivalent `float8` can. Also, the names `interval`,

`time`, and `timestamp` can only be used in this fashion if they are double-quoted, because of syntactic conflicts. Therefore, the use of the function-like cast syntax leads to inconsistencies and should probably be avoided.

Note

The function-like syntax is in fact just a function call. When one of the two standard cast syntaxes is used to do a run-time conversion, it will internally invoke a registered function to perform the conversion. By convention, these conversion functions have the same name as their output type, and thus the “function-like syntax” is nothing more than a direct invocation of the underlying conversion function. Obviously, this is not something that a portable application should rely on. For further details see [CREATE CAST](#).

4.2.10. Collation Expressions

The `COLLATE` clause overrides the collation of an expression. It is appended to the expression it applies to:

```
expr COLLATE collation
```

where *collation* is a possibly schema-qualified identifier. The `COLLATE` clause binds tighter than operators; parentheses can be used when necessary.

If no collation is explicitly specified, the database system either derives a collation from the columns involved in the expression, or it defaults to the default collation of the database if no column is involved in the expression.

The two common uses of the `COLLATE` clause are overriding the sort order in an `ORDER BY` clause, for example:

```
SELECT a, b, c FROM tbl WHERE ... ORDER BY a COLLATE "C";
```

and overriding the collation of a function or operator call that has locale-sensitive results, for example:

```
SELECT * FROM tbl WHERE a > 'foo' COLLATE "C";
```

Note that in the latter case the `COLLATE` clause is attached to an input argument of the operator we wish to affect. It doesn't matter which argument of the operator or function call the `COLLATE` clause is attached to, because the collation that is applied by the operator or function is derived by considering all arguments, and an explicit `COLLATE` clause will override the collations of all other arguments. (Attaching non-matching `COLLATE` clauses to more than one argument, however, is an error. For more details see [Section 23.2](#).) Thus, this gives the same result as the previous example:

```
SELECT * FROM tbl WHERE a COLLATE "C" > 'foo';
```

But this is an error:

```
SELECT * FROM tbl WHERE (a > 'foo') COLLATE "C";
```

because it attempts to apply a collation to the result of the `>` operator, which is of the non-collatable data type `boolean`.

4.2.11. Scalar Subqueries

A scalar subquery is an ordinary `SELECT` query in parentheses that returns exactly one row with one column. (See [Chapter 7](#) for information about writing queries.) The `SELECT` query is executed and the single returned value is used in the surrounding value expression. It is an error to use a query that returns more than one row or more than one column as a scalar subquery. (But if, during a particular execution, the subquery returns no rows, there is no error; the scalar result is taken to be null.) The subquery can refer to variables from the surrounding query, which will act as constants during any one evaluation of the subquery. See also [Section 9.24](#) for other expressions involving subqueries.

For example, the following finds the largest city population in each state:

```
SELECT name, (SELECT max(pop) FROM cities WHERE cities.state = states.name)
FROM states;
```

4.2.12. Array Constructors

An array constructor is an expression that builds an array value using values for its member elements. A simple array constructor consists of the key word `ARRAY`, a left square bracket `[`, a list of expressions (separated by commas) for the array element values, and finally a right square bracket `]`. For example:

```
SELECT ARRAY[1,2,3+4];
      array
-----
{1,2,7}
(1 row)
```

By default, the array element type is the common type of the member expressions, determined using the same rules as for `UNION` or `CASE` constructs (see [Section 10.5](#)). You can override this by explicitly casting the array constructor to the desired type, for example:

```
SELECT ARRAY[1,2,22.7]::integer[];
      array
-----
{1,2,23}
(1 row)
```

This has the same effect as casting each expression to the array element type individually. For more on casting, see [Section 4.2.9](#).

Multidimensional array values can be built by nesting array constructors. In the inner constructors, the key word `ARRAY` can be omitted. For example, these produce the same result:

```
SELECT ARRAY[ARRAY[1,2], ARRAY[3,4]];
      array
-----
{{1,2},{3,4}}
(1 row)
```

```
SELECT ARRAY[[1,2],[3,4]];
      array
-----
{{1,2},{3,4}}
(1 row)
```

Since multidimensional arrays must be rectangular, inner constructors at the same level must produce sub-arrays of identical dimensions. Any cast applied to the outer `ARRAY` constructor propagates automatically to all the inner constructors.

Multidimensional array constructor elements can be anything yielding an array of the proper kind, not only a sub-`ARRAY` construct. For example:

```
CREATE TABLE arr(f1 int[], f2 int[]);

INSERT INTO arr VALUES (ARRAY[[1,2],[3,4]], ARRAY[[5,6],[7,8]]);

SELECT ARRAY[f1, f2, '{{9,10},{11,12}}'::int[]] FROM arr;
      array
-----
{{1,2},{3,4}},{5,6},{7,8}},{9,10},{11,12}}
(1 row)
```

You can construct an empty array, but since it's impossible to have an array with no type, you must explicitly cast your empty array to the desired type. For example:

```
SELECT ARRAY[]::integer[];
array
-----
{}
(1 row)
```

It is also possible to construct an array from the results of a subquery. In this form, the array constructor is written with the key word `ARRAY` followed by a parenthesized (not bracketed) subquery. For example:

```
SELECT ARRAY(SELECT oid FROM pg_proc WHERE proname LIKE 'bytea%');
array
-----
{2011,1954,1948,1952,1951,1244,1950,2005,1949,1953,2006,31,2412}
(1 row)
```

```
SELECT ARRAY(SELECT ARRAY[i, i*2] FROM generate_series(1,5) AS a(i));
array
-----
{{1,2},{2,4},{3,6},{4,8},{5,10}}
(1 row)
```

The subquery must return a single column. If the subquery's output column is of a non-array type, the resulting one-dimensional array will have an element for each row in the subquery result, with an element type matching that of the subquery's output column. If the subquery's output column is of an array type, the result will be an array of the same type but one higher dimension; in this case all the subquery rows must yield arrays of identical dimensionality, else the result would not be rectangular.

The subscripts of an array value built with `ARRAY` always begin with one. For more information about arrays, see [Section 8.15](#).

4.2.13. Row Constructors

A row constructor is an expression that builds a row value (also called a composite value) using values for its member fields. A row constructor consists of the key word `ROW`, a left parenthesis, zero or more expressions (separated by commas) for the row field values, and finally a right parenthesis. For example:

```
SELECT ROW(1,2.5,'this is a test');
```

The key word `ROW` is optional when there is more than one expression in the list.

A row constructor can include the syntax `rowvalue.*`, which will be expanded to a list of the elements of the row value, just as occurs when the `.*` syntax is used at the top level of a `SELECT` list (see [Section 8.16.5](#)). For example, if table `t` has columns `f1` and `f2`, these are the same:

```
SELECT ROW(t.*, 42) FROM t;
SELECT ROW(t.f1, t.f2, 42) FROM t;
```

Note

Before PostgreSQL 8.2, the `.*` syntax was not expanded in row constructors, so that writing `ROW(t.*, 42)` created a two-field row whose first field was another row value. The new behavior is usually more useful. If you need the old behavior of nested row values, write the inner row value without `.*`, for instance `ROW(t, 42)`.

By default, the value created by a `ROW` expression is of an anonymous record type. If necessary, it can be cast to a named composite type — either the row type of a table, or a composite type created with `CREATE TYPE AS`. An explicit cast might be needed to avoid ambiguity. For example:

```
CREATE TABLE mytable(f1 int, f2 float, f3 text);
```

```
CREATE FUNCTION getf1(mytable) RETURNS int AS 'SELECT $1.f1' LANGUAGE SQL;

-- No cast needed since only one getf1() exists
SELECT getf1(ROW(1,2.5,'this is a test'));
  getf1
-----
      1
(1 row)

CREATE TYPE myrowtype AS (f1 int, f2 text, f3 numeric);

CREATE FUNCTION getf1(myrowtype) RETURNS int AS 'SELECT $1.f1' LANGUAGE SQL;

-- Now we need a cast to indicate which function to call:
SELECT getf1(ROW(1,2.5,'this is a test'));
ERROR:  function getf1(record) is not unique

SELECT getf1(ROW(1,2.5,'this is a test')::mytable);
  getf1
-----
      1
(1 row)

SELECT getf1(CAST(ROW(11,'this is a test',2.5) AS myrowtype));
  getf1
-----
     11
(1 row)
```

Row constructors can be used to build composite values to be stored in a composite-type table column, or to be passed to a function that accepts a composite parameter. Also, it is possible to test rows using the standard comparison operators as described in [Section 9.2](#), to compare one row against another as described in [Section 9.25](#), and to use them in connection with subqueries, as discussed in [Section 9.24](#),

4.2.14. Expression Evaluation Rules

The order of evaluation of subexpressions is not defined. In particular, the inputs of an operator or function are not necessarily evaluated left-to-right or in any other fixed order.

Furthermore, if the result of an expression can be determined by evaluating only some parts of it, then other subexpressions might not be evaluated at all. For instance, if one wrote:

```
SELECT true OR somefunc();
```

then `somefunc()` would (probably) not be called at all. The same would be the case if one wrote:

```
SELECT somefunc() OR true;
```

Note that this is not the same as the left-to-right “short-circuiting” of Boolean operators that is found in some programming languages.

As a consequence, it is unwise to use functions with side effects as part of complex expressions. It is particularly dangerous to rely on side effects or evaluation order in `WHERE` and `HAVING` clauses, since those clauses are extensively reprocessed as part of developing an execution plan. Boolean expressions (`AND/OR/NOT` combinations) in those clauses can be reorganized in any manner allowed by the laws of Boolean algebra.

When it is essential to force evaluation order, a `CASE` construct (see [Section 9.18](#)) can be used. For example, this is an untrustworthy way of trying to avoid division by zero in a `WHERE` clause:

```
SELECT ... WHERE x > 0 AND y/x > 1.5;
```

But this is safe:

```
SELECT ... WHERE CASE WHEN x > 0 THEN y/x > 1.5 ELSE false END;
```

A `CASE` construct used in this fashion will defeat optimization attempts, so it should only be done when necessary. (In this particular example, it would be better to sidestep the problem by writing `y > 1.5*x` instead.)

`CASE` is not a cure-all for such issues, however. One limitation of the technique illustrated above is that it does not prevent early evaluation of constant subexpressions. As described in [Section 39.7](#), functions and operators marked `IMMUTABLE` can be evaluated when the query is planned rather than when it is executed. Thus for example

```
SELECT CASE WHEN x > 0 THEN x ELSE 1/0 END FROM tab;
```

is likely to result in a division-by-zero failure due to the planner trying to simplify the constant subexpression, even if every row in the table has `x > 0` so that the `ELSE` arm would never be entered at run time.

While that particular example might seem silly, related cases that don't obviously involve constants can occur in queries executed within functions, since the values of function arguments and local variables can be inserted into queries as constants for planning purposes. Within PL/pgSQL functions, for example, using an `IF-THEN-ELSE` statement to protect a risky computation is much safer than just nesting it in a `CASE` expression.

Another limitation of the same kind is that a `CASE` cannot prevent evaluation of an aggregate expression contained within it, because aggregate expressions are computed before other expressions in a `SELECT` list or `HAVING` clause are considered. For example, the following query can cause a division-by-zero error despite seemingly having protected against it:

```
SELECT CASE WHEN min(employees) > 0
           THEN avg(expenses / employees)
           END
FROM departments;
```

The `min()` and `avg()` aggregates are computed concurrently over all the input rows, so if any row has `employees` equal to zero, the division-by-zero error will occur before there is any opportunity to test the result of `min()`. Instead, use a `WHERE` or `FILTER` clause to prevent problematic input rows from reaching an aggregate function in the first place.

4.3. Calling Functions

Postgres Pro allows functions that have named parameters to be called using either *positional* or *named* notation. Named notation is especially useful for functions that have a large number of parameters, since it makes the associations between parameters and actual arguments more explicit and reliable. In positional notation, a function call is written with its argument values in the same order as they are defined in the function declaration. In named notation, the arguments are matched to the function parameters by name and can be written in any order. For each notation, also consider the effect of function argument types, documented in [Section 10.3](#).

In either notation, parameters that have default values given in the function declaration need not be written in the call at all. But this is particularly useful in named notation, since any combination of parameters can be omitted; while in positional notation parameters can only be omitted from right to left.

Postgres Pro also supports *mixed* notation, which combines positional and named notation. In this case, positional parameters are written first and named parameters appear after them.

The following examples will illustrate the usage of all three notations, using the following function definition:

```
CREATE FUNCTION concat_lower_or_upper(a text, b text, uppercase boolean DEFAULT false)
RETURNS text
AS
```

```
$$
SELECT CASE
    WHEN $3 THEN UPPER($1 || ' ' || $2)
    ELSE LOWER($1 || ' ' || $2)
END;
$$
LANGUAGE SQL IMMUTABLE STRICT;
```

Function `concat_lower_or_upper` has two mandatory parameters, `a` and `b`. Additionally there is one optional parameter `uppercase` which defaults to `false`. The `a` and `b` inputs will be concatenated, and forced to either upper or lower case depending on the `uppercase` parameter. The remaining details of this function definition are not important here (see [Chapter 39](#) for more information).

4.3.1. Using Positional Notation

Positional notation is the traditional mechanism for passing arguments to functions in Postgres Pro. An example is:

```
SELECT concat_lower_or_upper('Hello', 'World', true);
concat_lower_or_upper
-----
HELLO WORLD
(1 row)
```

All arguments are specified in order. The result is upper case since `uppercase` is specified as `true`. Another example is:

```
SELECT concat_lower_or_upper('Hello', 'World');
concat_lower_or_upper
-----
hello world
(1 row)
```

Here, the `uppercase` parameter is omitted, so it receives its default value of `false`, resulting in lower case output. In positional notation, arguments can be omitted from right to left so long as they have defaults.

4.3.2. Using Named Notation

In named notation, each argument's name is specified using `=>` to separate it from the argument expression. For example:

```
SELECT concat_lower_or_upper(a => 'Hello', b => 'World');
concat_lower_or_upper
-----
hello world
(1 row)
```

Again, the argument `uppercase` was omitted so it is set to `false` implicitly. One advantage of using named notation is that the arguments may be specified in any order, for example:

```
SELECT concat_lower_or_upper(a => 'Hello', b => 'World', uppercase => true);
concat_lower_or_upper
-----
HELLO WORLD
(1 row)

SELECT concat_lower_or_upper(a => 'Hello', uppercase => true, b => 'World');
concat_lower_or_upper
-----
HELLO WORLD
(1 row)
```

An older syntax based on ":"=" is supported for backward compatibility:

```
SELECT concat_lower_or_upper(a := 'Hello', uppercase := true, b := 'World');
concat_lower_or_upper
-----
HELLO WORLD
(1 row)
```

4.3.3. Using Mixed Notation

The mixed notation combines positional and named notation. However, as already mentioned, named arguments cannot precede positional arguments. For example:

```
SELECT concat_lower_or_upper('Hello', 'World', uppercase => true);
concat_lower_or_upper
-----
HELLO WORLD
(1 row)
```

In the above query, the arguments `a` and `b` are specified positionally, while `uppercase` is specified by name. In this example, that adds little except documentation. With a more complex function having numerous parameters that have default values, named or mixed notation can save a great deal of writing and reduce chances for error.

Note

Named and mixed call notations currently cannot be used when calling an aggregate function (but they do work when an aggregate function is used as a window function).

Chapter 5. Data Definition

This chapter covers how one creates the database structures that will hold one's data. In a relational database, the raw data is stored in tables, so the majority of this chapter is devoted to explaining how tables are created and modified and what features are available to control what data is stored in the tables. Subsequently, we discuss how tables can be organized into schemas, and how privileges can be assigned to tables. Finally, we will briefly look at other features that affect the data storage, such as inheritance, table partitioning, views, functions, and triggers.

5.1. Table Basics

A table in a relational database is much like a table on paper: It consists of rows and columns. The number and order of the columns is fixed, and each column has a name. The number of rows is variable — it reflects how much data is stored at a given moment. SQL does not make any guarantees about the order of the rows in a table. When a table is read, the rows will appear in an unspecified order, unless sorting is explicitly requested. This is covered in [Chapter 7](#). Furthermore, SQL does not assign unique identifiers to rows, so it is possible to have several completely identical rows in a table. This is a consequence of the mathematical model that underlies SQL but is usually not desirable. Later in this chapter we will see how to deal with this issue.

Each column has a data type. The data type constrains the set of possible values that can be assigned to a column and assigns semantics to the data stored in the column so that it can be used for computations. For instance, a column declared to be of a numerical type will not accept arbitrary text strings, and the data stored in such a column can be used for mathematical computations. By contrast, a column declared to be of a character string type will accept almost any kind of data but it does not lend itself to mathematical calculations, although other operations such as string concatenation are available.

Postgres Pro includes a sizable set of built-in data types that fit many applications. Users can also define their own data types. Most built-in data types have obvious names and semantics, so we defer a detailed explanation to [Chapter 8](#). Some of the frequently used data types are `integer` for whole numbers, `numeric` for possibly fractional numbers, `text` for character strings, `date` for dates, `time` for time-of-day values, and `timestamp` for values containing both date and time.

To create a table, you use the aptly named `CREATE TABLE` command. In this command you specify at least a name for the new table, the names of the columns and the data type of each column. For example:

```
CREATE TABLE my_first_table (  
    first_column text,  
    second_column integer  
);
```

This creates a table named `my_first_table` with two columns. The first column is named `first_column` and has a data type of `text`; the second column has the name `second_column` and the type `integer`. The table and column names follow the identifier syntax explained in [Section 4.1.1](#). The type names are usually also identifiers, but there are some exceptions. Note that the column list is comma-separated and surrounded by parentheses.

Of course, the previous example was heavily contrived. Normally, you would give names to your tables and columns that convey what kind of data they store. So let's look at a more realistic example:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric  
);
```

(The `numeric` type can store fractional components, as would be typical of monetary amounts.)

Tip

When you create many interrelated tables it is wise to choose a consistent naming pattern for the tables and columns. For instance, there is a choice of using singular or plural nouns for table names, both of which are favored by some theorist or other.

There is a limit on how many columns a table can contain. Depending on the column types, it is between 250 and 1600. However, defining a table with anywhere near this many columns is highly unusual and often a questionable design.

If you no longer need a table, you can remove it using the [DROP TABLE](#) command. For example:

```
DROP TABLE my_first_table;
DROP TABLE products;
```

Attempting to drop a table that does not exist is an error. Nevertheless, it is common in SQL script files to unconditionally try to drop each table before creating it, ignoring any error messages, so that the script works whether or not the table exists. (If you like, you can use the `DROP TABLE IF EXISTS` variant to avoid the error messages, but this is not standard SQL.)

If you need to modify a table that already exists, see [Section 5.7](#) later in this chapter.

With the tools discussed so far you can create fully functional tables. The remainder of this chapter is concerned with adding features to the table definition to ensure data integrity, security, or convenience. If you are eager to fill your tables with data now you can skip ahead to [Chapter 6](#) and read the rest of this chapter later.

5.2. Default Values

A column can be assigned a default value. When a new row is created and no values are specified for some of the columns, those columns will be filled with their respective default values. A data manipulation command can also request explicitly that a column be set to its default value, without having to know what that value is. (Details about data manipulation commands are in [Chapter 6](#).)

If no default value is declared explicitly, the default value is the null value. This usually makes sense because a null value can be considered to represent unknown data.

In a table definition, default values are listed after the column data type. For example:

```
CREATE TABLE products (
    product_no integer,
    name text,
    price numeric DEFAULT 9.99
);
```

The default value can be an expression, which will be evaluated whenever the default value is inserted (*not* when the table is created). A common example is for a `timestamp` column to have a default of `CURRENT_TIMESTAMP`, so that it gets set to the time of row insertion. Another common example is generating a “serial number” for each row. In Postgres Pro this is typically done by something like:

```
CREATE TABLE products (
    product_no integer DEFAULT nextval('products_product_no_seq'),
    ...
);
```

where the `nextval()` function supplies successive values from a *sequence object* (see [Section 9.17](#)). This arrangement is sufficiently common that there's a special shorthand for it:

```
CREATE TABLE products (
    product_no SERIAL,
    ...
);
```

```
);
```

The `SERIAL` shorthand is discussed further in [Section 8.1.4](#).

5.3. Identity Columns

An identity column is a special column that is generated automatically from an implicit sequence. It can be used to generate key values. For distributed system specifics, refer to [CREATE TABLE](#).

To create an identity column, use the `GENERATED ... AS IDENTITY` clause in `CREATE TABLE`, for example:

```
CREATE TABLE people (  
    id bigint GENERATED ALWAYS AS IDENTITY,  
    ...  
);
```

or alternatively

```
CREATE TABLE people (  
    id bigint GENERATED BY DEFAULT AS IDENTITY,  
    ...  
);
```

See [CREATE TABLE](#) for more details.

If an `INSERT` command is executed on the table with the identity column and no value is explicitly specified for the identity column, then a value generated by the implicit sequence is inserted. For example, with the above definitions and assuming additional appropriate columns, writing

```
INSERT INTO people (name, address) VALUES ('A', 'foo');  
INSERT INTO people (name, address) VALUES ('B', 'bar');
```

would generate values for the `id` column starting at 1 and result in the following table data:

id	name	address
1	A	foo
2	B	bar

Alternatively, the keyword `DEFAULT` can be specified in place of a value to explicitly request the sequence-generated value, like

```
INSERT INTO people (id, name, address) VALUES (DEFAULT, 'C', 'baz');
```

Similarly, the keyword `DEFAULT` can be used in `UPDATE` commands.

Thus, in many ways, an identity column behaves like a column with a default value.

The clauses `ALWAYS` and `BY DEFAULT` in the column definition determine how explicitly user-specified values are handled in `INSERT` and `UPDATE` commands. In an `INSERT` command, if `ALWAYS` is selected, a user-specified value is only accepted if the `INSERT` statement specifies `OVERRIDING SYSTEM VALUE`. If `BY DEFAULT` is selected, then the user-specified value takes precedence. Thus, using `BY DEFAULT` results in a behavior more similar to default values, where the default value can be overridden by an explicit value, whereas `ALWAYS` provides some more protection against accidentally inserting an explicit value.

The data type of an identity column must be one of the data types supported by sequences. (See [CREATE SEQUENCE](#).) The properties of the associated sequence may be specified when creating an identity column (see [CREATE TABLE](#)) or changed afterwards (see [ALTER TABLE](#)).

An identity column is automatically marked as `NOT NULL`. An identity column, however, does not guarantee uniqueness. (A sequence normally returns unique values, but a sequence could be reset, or values could be inserted manually into the identity column, as discussed above.) Uniqueness would need to be enforced using a `PRIMARY KEY` or `UNIQUE` constraint.

In table inheritance hierarchies, identity columns and their properties in a child table are independent of those in its parent tables. A child table does not inherit identity columns or their properties automatically

from the parent. During `INSERT` or `UPDATE`, a column is treated as an identity column if that column is an identity column in the table named in the statement, and the corresponding identity properties are applied.

Partitions inherit identity columns from the partitioned table. They cannot have their own identity columns. The properties of a given identity column are consistent across all the partitions in the partition hierarchy.

5.4. Generated Columns

A generated column is a special column that is always computed from other columns. Thus, it is for columns what a view is for tables. There are two kinds of generated columns: stored and virtual. A stored generated column is computed when it is written (inserted or updated) and occupies storage as if it were a normal column. A virtual generated column occupies no storage and is computed when it is read. Thus, a virtual generated column is similar to a view and a stored generated column is similar to a materialized view (except that it is always updated automatically). Postgres Pro currently implements only stored generated columns.

For distributed system specifics, refer to [CREATE TABLE](#).

To create a generated column, use the `GENERATED ALWAYS AS` clause in `CREATE TABLE`, for example:

```
CREATE TABLE people (  
    ...,  
    height_cm numeric,  
    height_in numeric GENERATED ALWAYS AS (height_cm / 2.54) STORED  
);
```

The keyword `STORED` must be specified to choose the stored kind of generated column. See [CREATE TABLE](#) for more details.

A generated column cannot be written to directly. In `INSERT` or `UPDATE` commands, a value cannot be specified for a generated column, but the keyword `DEFAULT` may be specified.

Consider the differences between a column with a default and a generated column. The column default is evaluated once when the row is first inserted if no other value was provided; a generated column is updated whenever the row changes and cannot be overridden. A column default may not refer to other columns of the table; a generation expression would normally do so. A column default can use volatile functions, for example `random()` or functions referring to the current time; this is not allowed for generated columns.

Several restrictions apply to the definition of generated columns and tables involving generated columns:

- The generation expression can only use immutable functions and cannot use subqueries or reference anything other than the current row in any way.
- A generation expression cannot reference another generated column.
- A generation expression cannot reference a system column, except `tableoid`.
- A generated column cannot have a column default or an identity definition.
- A generated column cannot be part of a partition key.
- Foreign tables can have generated columns.
- For inheritance and partitioning:
 - If a parent column is a generated column, its child column must also be a generated column; however, the child column can have a different generation expression. The generation expression that is actually applied during insert or update of a row is the one associated with the table that the row is physically in. (This is unlike the behavior for column defaults: for those, the default value associated with the table named in the query applies.)
 - If a parent column is not a generated column, its child column must not be generated either.

- For inherited tables, if you write a child column definition without any `GENERATED` clause in `CREATE TABLE ... INHERITS`, then its `GENERATED` clause will automatically be copied from the parent. `ALTER TABLE ... INHERIT` will insist that parent and child columns already match as to generation status, but it will not require their generation expressions to match.
- Similarly for partitioned tables, if you write a child column definition without any `GENERATED` clause in `CREATE TABLE ... PARTITION OF`, then its `GENERATED` clause will automatically be copied from the parent. `ALTER TABLE ... ATTACH PARTITION` will insist that parent and child columns already match as to generation status, but it will not require their generation expressions to match.
- In case of multiple inheritance, if one parent column is a generated column, then all parent columns must be generated columns. If they do not all have the same generation expression, then the desired expression for the child must be specified explicitly.

Additional considerations apply to the use of generated columns.

- Generated columns maintain access privileges separately from their underlying base columns. So, it is possible to arrange it so that a particular role can read from a generated column but not from the underlying base columns.
- Generated columns are, conceptually, updated after `BEFORE` triggers have run. Therefore, changes made to base columns in a `BEFORE` trigger will be reflected in generated columns. But conversely, it is not allowed to access generated columns in `BEFORE` triggers.
- Generated columns are skipped for logical replication and cannot be specified in a `CREATE PUBLICATION` column list.

5.5. Constraints

Data types are a way to limit the kind of data that can be stored in a table. For many applications, however, the constraint they provide is too coarse. For example, a column containing a product price should probably only accept positive values. But there is no standard data type that accepts only positive numbers. Another issue is that you might want to constrain column data with respect to other columns or rows. For example, in a table containing product information, there should be only one row for each product number.

To that end, SQL allows you to define constraints on columns and tables. Constraints give you as much control over the data in your tables as you wish. If a user attempts to store data in a column that would violate a constraint, an error is raised. This applies even if the value came from the default value definition.

5.5.1. Check Constraints

A check constraint is the most generic constraint type. It allows you to specify that the value in a certain column must satisfy a Boolean (truth-value) expression. For instance, to require positive product prices, you could use:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric CHECK (price > 0)  
);
```

As you see, the constraint definition comes after the data type, just like default value definitions. Default values and constraints can be listed in any order. A check constraint consists of the key word `CHECK` followed by an expression in parentheses. The check constraint expression should involve the column thus constrained, otherwise the constraint would not make too much sense.

You can also give the constraint a separate name. This clarifies error messages and allows you to refer to the constraint when you need to change it. The syntax is:

```
CREATE TABLE products (  

```

```
product_no integer,  
name text,  
price numeric CONSTRAINT positive_price CHECK (price > 0)  
);
```

So, to specify a named constraint, use the key word `CONSTRAINT` followed by an identifier followed by the constraint definition. (If you don't specify a constraint name in this way, the system chooses a name for you.)

A check constraint can also refer to several columns. Say you store a regular price and a discounted price, and you want to ensure that the discounted price is lower than the regular price:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric CHECK (price > 0),  
    discounted_price numeric CHECK (discounted_price > 0),  
    CHECK (price > discounted_price)  
);
```

The first two constraints should look familiar. The third one uses a new syntax. It is not attached to a particular column, instead it appears as a separate item in the comma-separated column list. Column definitions and these constraint definitions can be listed in mixed order.

We say that the first two constraints are column constraints, whereas the third one is a table constraint because it is written separately from any one column definition. Column constraints can also be written as table constraints, while the reverse is not necessarily possible, since a column constraint is supposed to refer to only the column it is attached to. (Postgres Pro doesn't enforce that rule, but you should follow it if you want your table definitions to work with other database systems.) The above example could also be written as:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric,  
    CHECK (price > 0),  
    discounted_price numeric,  
    CHECK (discounted_price > 0),  
    CHECK (price > discounted_price)  
);
```

or even:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric CHECK (price > 0),  
    discounted_price numeric,  
    CHECK (discounted_price > 0 AND price > discounted_price)  
);
```

It's a matter of taste.

Names can be assigned to table constraints in the same way as column constraints:

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric,  
    CHECK (price > 0),  
    discounted_price numeric,  
    CHECK (discounted_price > 0),  
    CONSTRAINT valid_discount CHECK (price > discounted_price)
```

```
);
```

It should be noted that a check constraint is satisfied if the check expression evaluates to true or the null value. Since most expressions will evaluate to the null value if any operand is null, they will not prevent null values in the constrained columns. To ensure that a column does not contain null values, the not-null constraint described in the next section can be used.

Note

Postgres Pro does not support `CHECK` constraints that reference table data other than the new or updated row being checked. While a `CHECK` constraint that violates this rule may appear to work in simple tests, it cannot guarantee that the database will not reach a state in which the constraint condition is false (due to subsequent changes of the other row(s) involved). This would cause a database dump and restore to fail. The restore could fail even when the complete database state is consistent with the constraint, due to rows not being loaded in an order that will satisfy the constraint. If possible, use `UNIQUE`, `EXCLUDE`, or `FOREIGN KEY` constraints to express cross-row and cross-table restrictions.

If what you desire is a one-time check against other rows at row insertion, rather than a continuously-maintained consistency guarantee, a custom [trigger](#) can be used to implement that. (This approach avoids the dump/restore problem because `pg_dump` does not reinstall triggers until after restoring data, so that the check will not be enforced during a dump/restore.)

Note

Postgres Pro assumes that `CHECK` constraints' conditions are immutable, that is, they will always give the same result for the same input row. This assumption is what justifies examining `CHECK` constraints only when rows are inserted or updated, and not at other times. (The warning above about not referencing other table data is really a special case of this restriction.)

An example of a common way to break this assumption is to reference a user-defined function in a `CHECK` expression, and then change the behavior of that function. Postgres Pro does not disallow that, but it will not notice if there are rows in the table that now violate the `CHECK` constraint. That would cause a subsequent database dump and restore to fail. The recommended way to handle such a change is to drop the constraint (using `ALTER TABLE`), adjust the function definition, and re-add the constraint, thereby rechecking it against all table rows.

5.5.2. Not-Null Constraints

A not-null constraint simply specifies that a column must not assume the null value. A syntax example:

```
CREATE TABLE products (  
    product_no integer NOT NULL,  
    name text NOT NULL,  
    price numeric  
);
```

A not-null constraint is always written as a column constraint. A not-null constraint is functionally equivalent to creating a check constraint `CHECK (column_name IS NOT NULL)`, but in Postgres Pro creating an explicit not-null constraint is more efficient. The drawback is that you cannot give explicit names to not-null constraints created this way.

Of course, a column can have more than one constraint. Just write the constraints one after another:

```
CREATE TABLE products (  
    product_no integer NOT NULL,  
    name text NOT NULL,
```

```
    price numeric NOT NULL CHECK (price > 0)
);
```

The order doesn't matter. It does not necessarily determine in which order the constraints are checked.

The `NOT NULL` constraint has an inverse: the `NULL` constraint. This does not mean that the column must be null, which would surely be useless. Instead, this simply selects the default behavior that the column might be null. The `NULL` constraint is not present in the SQL standard and should not be used in portable applications. (It was only added to Postgres Pro to be compatible with some other database systems.) Some users, however, like it because it makes it easy to toggle the constraint in a script file. For example, you could start with:

```
CREATE TABLE products (
    product_no integer NULL,
    name text NULL,
    price numeric NULL
);
```

and then insert the `NOT` key word where desired.

Tip

In most database designs the majority of columns should be marked not null.

5.5.3. Unique Constraints

Unique constraints ensure that the data contained in a column, or a group of columns, is unique among all the rows in the table. The syntax is:

```
CREATE TABLE products (
    product_no integer UNIQUE,
    name text,
    price numeric
);
```

when written as a column constraint, and:

```
CREATE TABLE products (
    product_no integer,
    name text,
    price numeric,
    UNIQUE (product_no)
);
```

when written as a table constraint.

To define a unique constraint for a group of columns, write it as a table constraint with the column names separated by commas:

```
CREATE TABLE example (
    a integer,
    b integer,
    c integer,
    UNIQUE (a, c)
);
```

This specifies that the combination of values in the indicated columns is unique across the whole table, though any one of the columns need not be (and ordinarily isn't) unique.

You can assign your own name for a unique constraint, in the usual way:

```
CREATE TABLE products (
    product_no integer CONSTRAINT must_be_different UNIQUE,
```

```
    name text,  
    price numeric  
);
```

Adding a unique constraint will automatically create a unique B-tree index on the column or group of columns listed in the constraint. A uniqueness restriction covering only some rows cannot be written as a unique constraint, but it is possible to enforce such a restriction by creating a unique [partial index](#).

In general, a unique constraint is violated if there is more than one row in the table where the values of all of the columns included in the constraint are equal. By default, two null values are not considered equal in this comparison. That means even in the presence of a unique constraint it is possible to store duplicate rows that contain a null value in at least one of the constrained columns. This behavior can be changed by adding the clause `NULLS NOT DISTINCT`, like

```
CREATE TABLE products (  
    product_no integer UNIQUE NULLS NOT DISTINCT,  
    name text,  
    price numeric  
);  
  
or
```

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric,  
    UNIQUE NULLS NOT DISTINCT (product_no)  
);
```

The default behavior can be specified explicitly using `NULLS DISTINCT`. The default null treatment in unique constraints is implementation-defined according to the SQL standard, and other implementations have a different behavior. So be careful when developing applications that are intended to be portable.

5.5.4. Primary Keys

A primary key constraint indicates that a column, or group of columns, can be used as a unique identifier for rows in the table. This requires that the values be both unique and not null. So, the following two table definitions accept the same data:

```
CREATE TABLE products (  
    product_no integer UNIQUE NOT NULL,  
    name text,  
    price numeric  
);  
  
CREATE TABLE products (  
    product_no integer PRIMARY KEY,  
    name text,  
    price numeric  
);
```

Primary keys can span more than one column; the syntax is similar to unique constraints:

```
CREATE TABLE example (  
    a integer,  
    b integer,  
    c integer,  
    PRIMARY KEY (a, c)  
);
```

Adding a primary key will automatically create a unique B-tree index on the column or group of columns listed in the primary key, and will force the column(s) to be marked `NOT NULL`.

A table can have at most one primary key. (There can be any number of unique and not-null constraints, which are functionally almost the same thing, but only one can be identified as the primary key.) Relational database theory dictates that every table must have a primary key. This rule is not enforced by Postgres Pro, but it is usually best to follow it.

Primary keys are useful both for documentation purposes and for client applications. For example, a GUI application that allows modifying row values probably needs to know the primary key of a table to be able to identify rows uniquely. There are also various ways in which the database system makes use of a primary key if one has been declared; for example, the primary key defines the default target column(s) for foreign keys referencing its table.

5.5.5. Foreign Keys

A foreign key constraint specifies that the values in a column (or a group of columns) must match the values appearing in some row of another table. We say this maintains the *referential integrity* between two related tables.

Say you have the product table that we have used several times already:

```
CREATE TABLE products (  
    product_no integer PRIMARY KEY,  
    name text,  
    price numeric  
);
```

Let's also assume you have a table storing orders of those products. We want to ensure that the orders table only contains orders of products that actually exist. So we define a foreign key constraint in the orders table that references the products table:

```
CREATE TABLE orders (  
    order_id integer PRIMARY KEY,  
    product_no integer REFERENCES products (product_no),  
    quantity integer  
);
```

Now it is impossible to create orders with non-NULL `product_no` entries that do not appear in the products table.

We say that in this situation the orders table is the *referencing* table and the products table is the *referenced* table. Similarly, there are referencing and referenced columns.

You can also shorten the above command to:

```
CREATE TABLE orders (  
    order_id integer PRIMARY KEY,  
    product_no integer REFERENCES products,  
    quantity integer  
);
```

because in absence of a column list the primary key of the referenced table is used as the referenced column(s).

You can assign your own name for a foreign key constraint, in the usual way.

A foreign key can also constrain and reference a group of columns. As usual, it then needs to be written in table constraint form. Here is a contrived syntax example:

```
CREATE TABLE t1 (  
    a integer PRIMARY KEY,  
    b integer,  
    c integer,  
    FOREIGN KEY (b, c) REFERENCES other_table (c1, c2)  
);
```

Of course, the number and type of the constrained columns need to match the number and type of the referenced columns.

Sometimes it is useful for the “other table” of a foreign key constraint to be the same table; this is called a *self-referential* foreign key. For example, if you want rows of a table to represent nodes of a tree structure, you could write

```
CREATE TABLE tree (  
    node_id integer PRIMARY KEY,  
    parent_id integer REFERENCES tree,  
    name text,  
    ...  
);
```

A top-level node would have NULL `parent_id`, while non-NULL `parent_id` entries would be constrained to reference valid rows of the table.

A table can have more than one foreign key constraint. This is used to implement many-to-many relationships between tables. Say you have tables about products and orders, but now you want to allow one order to contain possibly many products (which the structure above did not allow). You could use this table structure:

```
CREATE TABLE products (  
    product_no integer PRIMARY KEY,  
    name text,  
    price numeric  
);  
  
CREATE TABLE orders (  
    order_id integer PRIMARY KEY,  
    shipping_address text,  
    ...  
);  
  
CREATE TABLE order_items (  
    product_no integer REFERENCES products,  
    order_id integer REFERENCES orders,  
    quantity integer,  
    PRIMARY KEY (product_no, order_id)  
);
```

Notice that the primary key overlaps with the foreign keys in the last table.

We know that the foreign keys disallow creation of orders that do not relate to any products. But what if a product is removed after an order is created that references it? SQL allows you to handle that as well. Intuitively, we have a few options:

- Disallow deleting a referenced product
- Delete the orders as well
- Something else?

To illustrate this, let's implement the following policy on the many-to-many relationship example above: when someone wants to remove a product that is still referenced by an order (via `order_items`), we disallow it. If someone removes an order, the order items are removed as well:

```
CREATE TABLE products (  
    product_no integer PRIMARY KEY,  
    name text,  
    price numeric  
);
```

```
CREATE TABLE orders (  
    order_id integer PRIMARY KEY,  
    shipping_address text,  
    ...  
);  
  
CREATE TABLE order_items (  
    product_no integer REFERENCES products ON DELETE RESTRICT,  
    order_id integer REFERENCES orders ON DELETE CASCADE,  
    quantity integer,  
    PRIMARY KEY (product_no, order_id)  
);
```

Restricting and cascading deletes are the two most common options. `RESTRICT` prevents deletion of a referenced row. `NO ACTION` means that if any referencing rows still exist when the constraint is checked, an error is raised; this is the default behavior if you do not specify anything. (The essential difference between these two choices is that `NO ACTION` allows the check to be deferred until later in the transaction, whereas `RESTRICT` does not.) `CASCADE` specifies that when a referenced row is deleted, row(s) referencing it should be automatically deleted as well. There are two other options: `SET NULL` and `SET DEFAULT`. These cause the referencing column(s) in the referencing row(s) to be set to nulls or their default values, respectively, when the referenced row is deleted. Note that these do not excuse you from observing any constraints. For example, if an action specifies `SET DEFAULT` but the default value would not satisfy the foreign key constraint, the operation will fail.

The appropriate choice of `ON DELETE` action depends on what kinds of objects the related tables represent. When the referencing table represents something that is a component of what is represented by the referenced table and cannot exist independently, then `CASCADE` could be appropriate. If the two tables represent independent objects, then `RESTRICT` or `NO ACTION` is more appropriate; an application that actually wants to delete both objects would then have to be explicit about this and run two delete commands. In the above example, order items are part of an order, and it is convenient if they are deleted automatically if an order is deleted. But products and orders are different things, and so making a deletion of a product automatically cause the deletion of some order items could be considered problematic. The actions `SET NULL` or `SET DEFAULT` can be appropriate if a foreign-key relationship represents optional information. For example, if the products table contained a reference to a product manager, and the product manager entry gets deleted, then setting the product's product manager to null or a default might be useful.

The actions `SET NULL` and `SET DEFAULT` can take a column list to specify which columns to set. Normally, all columns of the foreign-key constraint are set; setting only a subset is useful in some special cases. Consider the following example:

```
CREATE TABLE tenants (  
    tenant_id integer PRIMARY KEY  
);  
  
CREATE TABLE users (  
    tenant_id integer REFERENCES tenants ON DELETE CASCADE,  
    user_id integer NOT NULL,  
    PRIMARY KEY (tenant_id, user_id)  
);  
  
CREATE TABLE posts (  
    tenant_id integer REFERENCES tenants ON DELETE CASCADE,  
    post_id integer NOT NULL,  
    author_id integer,  
    PRIMARY KEY (tenant_id, post_id),  
    FOREIGN KEY (tenant_id, author_id) REFERENCES users ON DELETE SET NULL (author_id)  
);
```

Without the specification of the column, the foreign key would also set the column `tenant_id` to null, but that column is still required as part of the primary key.

Analogous to `ON DELETE` there is also `ON UPDATE` which is invoked when a referenced column is changed (updated). The possible actions are the same, except that column lists cannot be specified for `SET NULL` and `SET DEFAULT`. In this case, `CASCADE` means that the updated values of the referenced column(s) should be copied into the referencing row(s).

Normally, a referencing row need not satisfy the foreign key constraint if any of its referencing columns are null. If `MATCH FULL` is added to the foreign key declaration, a referencing row escapes satisfying the constraint only if all its referencing columns are null (so a mix of null and non-null values is guaranteed to fail a `MATCH FULL` constraint). If you don't want referencing rows to be able to avoid satisfying the foreign key constraint, declare the referencing column(s) as `NOT NULL`.

A foreign key must reference columns that either are a primary key or form a unique constraint, or are columns from a non-partial unique index. This means that the referenced columns always have an index to allow efficient lookups on whether a referencing row has a match. Since a `DELETE` of a row from the referenced table or an `UPDATE` of a referenced column will require a scan of the referencing table for rows matching the old value, it is often a good idea to index the referencing columns too. Because this is not always needed, and there are many choices available on how to index, the declaration of a foreign key constraint does not automatically create an index on the referencing columns.

More information about updating and deleting data is in [Chapter 6](#). Also see the description of foreign key constraint syntax in the reference documentation for [CREATE TABLE](#).

5.5.6. Exclusion Constraints

Exclusion constraints ensure that if any two rows are compared on the specified columns or expressions using the specified operators, at least one of these operator comparisons will return false or null. The syntax is:

```
CREATE TABLE circles (  
    c circle,  
    EXCLUDE USING gist (c WITH &&)  
);
```

See also [CREATE TABLE ... CONSTRAINT ... EXCLUDE](#) for details.

Adding an exclusion constraint will automatically create an index of the type specified in the constraint declaration.

5.6. System Columns

Every table has several *system columns* that are implicitly defined by the system. Therefore, these names cannot be used as names of user-defined columns. (Note that these restrictions are separate from whether the name is a key word or not; quoting a name will not allow you to escape these restrictions.) You do not really need to be concerned about these columns; just know they exist.

`tableoid`

The OID of the table containing this row. This column is particularly handy for queries that select from partitioned tables (see [Section 5.12](#)) or inheritance hierarchies (see [Section 5.11](#)), since without it, it's difficult to tell which individual table a row came from. The `tableoid` can be joined against the `oid` column of `pg_class` to obtain the table name.

`xmin`

The identity (transaction ID) of the inserting transaction for this row version. (A row version is an individual state of a row; each update of a row creates a new row version for the same logical row.)

`cmin`

The command identifier (starting at zero) within the inserting transaction.

`xmax`

The identity (transaction ID) of the deleting transaction, or zero for an undeleted row version. It is possible for this column to be nonzero in a visible row version. That usually indicates that the deleting transaction hasn't committed yet, or that an attempted deletion was rolled back.

`cmax`

The command identifier within the deleting transaction, or zero.

`ctid`

The physical location of the row version within its table. Note that although the `ctid` can be used to locate the row version very quickly, a row's `ctid` will change if it is updated or moved by `VACUUM FULL`. Therefore `ctid` is useless as a long-term row identifier. A primary key should be used to identify logical rows.

In Postgres Pro Shardman, transaction IDs are implemented as 64-bit counters to prevent transaction ID wraparound. For details, see [Section 24.1.5](#).

Command identifiers are 32-bit quantities. This creates a hard limit of 2^{32} (4 billion) SQL commands within a single transaction. In practice this limit is not a problem — note that the limit is on the number of SQL commands, not the number of rows processed. Also, only commands that actually modify the database contents will consume a command identifier.

5.7. Modifying Tables

When you create a table and you realize that you made a mistake, or the requirements of the application change, you can drop the table and create it again. But this is not a convenient option if the table is already filled with data, or if the table is referenced by other database objects (for instance a foreign key constraint). Therefore Postgres Pro provides a family of commands to make modifications to existing tables. Note that this is conceptually distinct from altering the data contained in the table: here we are interested in altering the definition, or structure, of the table.

You can:

- Add columns
- Remove columns
- Add constraints
- Remove constraints
- Change default values
- Change column data types
- Rename columns
- Rename tables

All these actions are performed using the [ALTER TABLE](#) command, whose reference page contains details beyond those given here.

5.7.1. Adding a Column

To add a column, use a command like:

```
ALTER TABLE products ADD COLUMN description text;
```

The new column is initially filled with whatever default value is given (null if you don't specify a `DEFAULT` clause).

Tip

From Postgres Pro 11, adding a column with a constant default value no longer means that each row of the table needs to be updated when the `ALTER TABLE` statement is executed. Instead, the default value will be returned the next time the row is accessed, and applied when the table is rewritten, making the `ALTER TABLE` very fast even on large tables.

However, if the default value is volatile (e.g., `clock_timestamp()`) each row will need to be updated with the value calculated at the time `ALTER TABLE` is executed. To avoid a potentially lengthy update operation, particularly if you intend to fill the column with mostly nondefault values anyway, it may be preferable to add the column with no default, insert the correct values using `UPDATE`, and then add any desired default as described below.

You can also define constraints on the column at the same time, using the usual syntax:

```
ALTER TABLE products ADD COLUMN description text CHECK (description <> '');
```

In fact all the options that can be applied to a column description in `CREATE TABLE` can be used here. Keep in mind however that the default value must satisfy the given constraints, or the `ADD` will fail. Alternatively, you can add constraints later (see below) after you've filled in the new column correctly.

5.7.2. Removing a Column

To remove a column, use a command like:

```
ALTER TABLE products DROP COLUMN description;
```

Whatever data was in the column disappears. Table constraints involving the column are dropped, too. However, if the column is referenced by a foreign key constraint of another table, Postgres Pro will not silently drop that constraint. You can authorize dropping everything that depends on the column by adding `CASCADE`:

```
ALTER TABLE products DROP COLUMN description CASCADE;
```

See [Section 5.15](#) for a description of the general mechanism behind this.

5.7.3. Adding a Constraint

To add a constraint, the table constraint syntax is used. For example:

```
ALTER TABLE products ADD CHECK (name <> '');  
ALTER TABLE products ADD CONSTRAINT some_name UNIQUE (product_no);  
ALTER TABLE products ADD FOREIGN KEY (product_group_id) REFERENCES product_groups;
```

To add a not-null constraint, which cannot be written as a table constraint, use this syntax:

```
ALTER TABLE products ALTER COLUMN product_no SET NOT NULL;
```

The constraint will be checked immediately, so the table data must satisfy the constraint before it can be added.

5.7.4. Removing a Constraint

To remove a constraint you need to know its name. If you gave it a name then that's easy. Otherwise the system assigned a generated name, which you need to find out. The `psql` command `\d tablename` can be helpful here; other interfaces might also provide a way to inspect table details. Then the command is:

```
ALTER TABLE products DROP CONSTRAINT some_name;
```

As with dropping a column, you need to add `CASCADE` if you want to drop a constraint that something else depends on. An example is that a foreign key constraint depends on a unique or primary key constraint on the referenced column(s).

This works the same for all constraint types except not-null constraints. To drop a not-null constraint use:

```
ALTER TABLE products ALTER COLUMN product_no DROP NOT NULL;
```

(Recall that not-null constraints do not have names.)

5.7.5. Changing a Column's Default Value

To set a new default for a column, use a command like:

```
ALTER TABLE products ALTER COLUMN price SET DEFAULT 7.77;
```

Note that this doesn't affect any existing rows in the table, it just changes the default for future `INSERT` commands.

To remove any default value, use:

```
ALTER TABLE products ALTER COLUMN price DROP DEFAULT;
```

This is effectively the same as setting the default to null. As a consequence, it is not an error to drop a default where one hadn't been defined, because the default is implicitly the null value.

5.7.6. Changing a Column's Data Type

To convert a column to a different data type, use a command like:

```
ALTER TABLE products ALTER COLUMN price TYPE numeric(10,2);
```

This will succeed only if each existing entry in the column can be converted to the new type by an implicit cast. If a more complex conversion is needed, you can add a `USING` clause that specifies how to compute the new values from the old.

Postgres Pro will attempt to convert the column's default value (if any) to the new type, as well as any constraints that involve the column. But these conversions might fail, or might produce surprising results. It's often best to drop any constraints on the column before altering its type, and then add back suitably modified constraints afterwards.

5.7.7. Renaming a Column

To rename a column:

```
ALTER TABLE products RENAME COLUMN product_no TO product_number;
```

5.7.8. Renaming a Table

To rename a table:

```
ALTER TABLE products RENAME TO items;
```

5.8. Privileges

When an object is created, it is assigned an owner. The owner is normally the role that executed the creation statement. For most kinds of objects, the initial state is that only the owner (or a superuser) can do anything with the object. To allow other roles to use it, *privileges* must be granted.

There are different kinds of privileges: `SELECT`, `INSERT`, `UPDATE`, `DELETE`, `TRUNCATE`, `REFERENCES`, `TRIGGER`, `CREATE`, `CONNECT`, `TEMPORARY`, `EXECUTE`, `USAGE`, `SET`, `ALTER SYSTEM`, and `MAINTAIN`. The privileges applicable to a particular object vary depending on the object's type (table, function, etc.). More detail about the meanings of these privileges appears below. The following sections and chapters will also show you how these privileges are used.

The right to modify or destroy an object is inherent in being the object's owner, and cannot be granted or revoked in itself. (However, like all privileges, that right can be inherited by members of the owning role; see [Section 21.3](#).)

An object can be assigned to a new owner with an `ALTER` command of the appropriate kind for the object, for example

```
ALTER TABLE table_name OWNER TO new_owner;
```

Superusers can always do this; ordinary roles can only do it if they are both the current owner of the object (or inherit the privileges of the owning role) and able to `SET ROLE` to the new owning role.

To assign privileges, the `GRANT` command is used. For example, if `joe` is an existing role, and `accounts` is an existing table, the privilege to update the table can be granted with:

```
GRANT UPDATE ON accounts TO joe;
```

Writing `ALL` in place of a specific privilege grants all privileges that are relevant for the object type.

The special “role” name `PUBLIC` can be used to grant a privilege to every role on the system. Also, “group” roles can be set up to help manage privileges when there are many users of a database — for details see [Chapter 21](#).

To revoke a previously-granted privilege, use the fittingly named `REVOKE` command:

```
REVOKE ALL ON accounts FROM PUBLIC;
```

Ordinarily, only the object's owner (or a superuser) can grant or revoke privileges on an object. However, it is possible to grant a privilege “with grant option”, which gives the recipient the right to grant it in turn to others. If the grant option is subsequently revoked then all who received the privilege from that recipient (directly or through a chain of grants) will lose the privilege. For details see the [GRANT](#) and [REVOKE](#) reference pages.

An object's owner can choose to revoke their own ordinary privileges, for example to make a table read-only for themselves as well as others. But owners are always treated as holding all grant options, so they can always re-grant their own privileges.

The available privileges are:

`SELECT`

Allows `SELECT` from any column, or specific column(s), of a table, view, materialized view, or other table-like object. Also allows use of `COPY TO`. This privilege is also needed to reference existing column values in `UPDATE`, `DELETE`, or `MERGE`. For sequences, this privilege also allows use of the `currval` function. For large objects, this privilege allows the object to be read.

`INSERT`

Allows `INSERT` of a new row into a table, view, etc. Can be granted on specific column(s), in which case only those columns may be assigned to in the `INSERT` command (other columns will therefore receive default values). Also allows use of `COPY FROM`.

`UPDATE`

Allows `UPDATE` of any column, or specific column(s), of a table, view, etc. (In practice, any nontrivial `UPDATE` command will require `SELECT` privilege as well, since it must reference table columns to determine which rows to update, and/or to compute new values for columns.) `SELECT ... FOR UPDATE` and `SELECT ... FOR SHARE` also require this privilege on at least one column, in addition to the `SELECT` privilege. For sequences, this privilege allows use of the `nextval` and `setval` functions. For large objects, this privilege allows writing or truncating the object.

`DELETE`

Allows `DELETE` of a row from a table, view, etc. (In practice, any nontrivial `DELETE` command will require `SELECT` privilege as well, since it must reference table columns to determine which rows to delete.)

`TRUNCATE`

Allows `TRUNCATE` on a table.

`REFERENCES`

Allows creation of a foreign key constraint referencing a table, or specific column(s) of a table.

`TRIGGER`

Allows creation of a trigger on a table, view, etc.

`CREATE`

For databases, allows new schemas and publications to be created within the database, and allows trusted extensions to be installed within the database.

For schemas, allows new objects to be created within the schema. To rename an existing object, you must own the object *and* have this privilege for the containing schema.

For tablespaces, allows tables, indexes, and temporary files to be created within the tablespace, and allows databases to be created that have the tablespace as their default tablespace.

Note that revoking this privilege will not alter the existence or location of existing objects.

CONNECT

Allows the grantee to connect to the database. This privilege is checked at connection startup (in addition to checking any restrictions imposed by `pg_hba.conf`).

TEMPORARY

Allows temporary tables to be created while using the database.

EXECUTE

Allows calling a function or procedure, including use of any operators that are implemented on top of the function. This is the only type of privilege that is applicable to functions and procedures.

USAGE

For procedural languages, allows use of the language for the creation of functions in that language. This is the only type of privilege that is applicable to procedural languages.

For schemas, allows access to objects contained in the schema (assuming that the objects' own privilege requirements are also met). Essentially this allows the grantee to “look up” objects within the schema. Without this permission, it is still possible to see the object names, e.g., by querying system catalogs. Also, after revoking this permission, existing sessions might have statements that have previously performed this lookup, so this is not a completely secure way to prevent object access.

For sequences, allows use of the `currval` and `nextval` functions.

For types and domains, allows use of the type or domain in the creation of tables, functions, and other schema objects. (Note that this privilege does not control all “usage” of the type, such as values of the type appearing in queries. It only prevents objects from being created that depend on the type. The main purpose of this privilege is controlling which users can create dependencies on a type, which could prevent the owner from changing the type later.)

For foreign-data wrappers, allows creation of new servers using the foreign-data wrapper.

For foreign servers, allows creation of foreign tables using the server. Grantees may also create, alter, or drop their own user mappings associated with that server.

SET

Allows a server configuration parameter to be set to a new value within the current session. (While this privilege can be granted on any parameter, it is meaningless except for parameters that would normally require superuser privilege to set.)

MAINTAIN

Allows `VACUUM`, `ANALYZE`, `CLUSTER`, `REFRESH MATERIALIZED VIEW`, `REINDEX`, and `LOCK TABLE` on a relation.

The privileges required by other commands are listed on the reference page of the respective command.

Postgres Pro grants privileges on some types of objects to `PUBLIC` by default when the objects are created. No privileges are granted to `PUBLIC` by default on tables, table columns, sequences, foreign data wrappers, foreign servers, large objects, schemas, tablespaces, or configuration parameters. For other types of objects, the default privileges granted to `PUBLIC` are as follows: `CONNECT` and `TEMPORARY` (create temporary tables) privileges for databases; `EXECUTE` privilege for functions and procedures; and `USAGE`

privilege for languages and data types (including domains). The object owner can, of course, `REVOKE` both default and expressly granted privileges. (For maximum security, issue the `REVOKE` in the same transaction that creates the object; then there is no window in which another user can use the object.) Also, these default privilege settings can be overridden using the `ALTER DEFAULT PRIVILEGES` command.

Table 5.1 shows the one-letter abbreviations that are used for these privilege types in *ACL* (Access Control List) values. You will see these letters in the output of the `psql` commands listed below, or when looking at ACL columns of system catalogs.

Table 5.1. ACL Privilege Abbreviations

Privilege	Abbreviation	Applicable Object Types
SELECT	r (“read”)	LARGE OBJECT, SEQUENCE, TABLE (and table-like objects), table column
INSERT	a (“append”)	TABLE, table column
UPDATE	w (“write”)	LARGE OBJECT, SEQUENCE, TABLE, table column
DELETE	d	TABLE
TRUNCATE	D	TABLE
REFERENCES	x	TABLE, table column
TRIGGER	t	TABLE
CREATE	C	DATABASE, SCHEMA, TABLESPACE
CONNECT	c	DATABASE
TEMPORARY	T	DATABASE
EXECUTE	X	FUNCTION, PROCEDURE
USAGE	U	DOMAIN, FOREIGN DATA WRAPPER, FOREIGN SERVER, LANGUAGE, SCHEMA, SEQUENCE, TYPE
SET	s	PARAMETER
ALTER SYSTEM	A	PARAMETER
MAINTAIN	m	TABLE

Table 5.2 summarizes the privileges available for each type of SQL object, using the abbreviations shown above. It also shows the `psql` command that can be used to examine privilege settings for each object type.

Table 5.2. Summary of Access Privileges

Object Type	All Privileges	Default <small>PUBLIC</small> Privileges	<code>psql</code> Command
DATABASE	CTc	Tc	<code>\l</code>
DOMAIN	U	U	<code>\dD+</code>
FUNCTION or PROCEDURE	X	X	<code>\df+</code>
FOREIGN DATA WRAPPER	U	none	<code>\dew+</code>
FOREIGN SERVER	U	none	<code>\des+</code>
LANGUAGE	U	U	<code>\dL+</code>
LARGE OBJECT	rw	none	<code>\dl+</code>
PARAMETER	sA	none	<code>\dconfig+</code>
SCHEMA	UC	none	<code>\dn+</code>
SEQUENCE	rwU	none	<code>\dp</code>
TABLE (and table-like objects)	arwdDxtm	none	<code>\dp</code>

Object Type	All Privileges	Default <small>PUBLIC</small> Privileges	psql Command
Table column	arwx	none	\dp
TABLESPACE	C	none	\db+
TYPE	U	U	\dT+

The privileges that have been granted for a particular object are displayed as a list of `aclitem` entries, each having the format:

```
grantee=privilege-abbreviation[*].../grantor
```

Each `aclitem` lists all the permissions of one grantee that have been granted by a particular grantor. Specific privileges are represented by one-letter abbreviations from [Table 5.1](#), with `*` appended if the privilege was granted with grant option. For example, `calvin=r*w/hobbes` specifies that the role `calvin` has the privilege `SELECT (r)` with grant option (`*`) as well as the non-grantable privilege `UPDATE (w)`, both granted by the role `hobbes`. If `calvin` also has some privileges on the same object granted by a different grantor, those would appear as a separate `aclitem` entry. An empty grantee field in an `aclitem` stands for `PUBLIC`.

As an example, suppose that user `miriam` creates table `mytable` and does:

```
GRANT SELECT ON mytable TO PUBLIC;
GRANT SELECT, UPDATE, INSERT ON mytable TO admin;
GRANT SELECT (col1), UPDATE (col1) ON mytable TO miriam_rw;
```

Then `psql`'s `\dp` command would show:

```
=> \dp mytable
```

Access privileges					
Schema	Name	Type	Access privileges	Column privileges	Policies
public	mytable	table	miriam=arwdDxtm/miriam+	col1:	+
			=r/miriam	miriam_rw=rw/miriam	
			admin=arw/miriam		

(1 row)

If the “Access privileges” column is empty for a given object, it means the object has default privileges (that is, its privileges entry in the relevant system catalog is null). Default privileges always include all privileges for the owner, and can include some privileges for `PUBLIC` depending on the object type, as explained above. The first `GRANT` or `REVOKE` on an object will instantiate the default privileges (producing, for example, `miriam=arwdDxt/miriam`) and then modify them per the specified request. Similarly, entries are shown in “Column privileges” only for columns with nondefault privileges. (Note: for this purpose, “default privileges” always means the built-in default privileges for the object's type. An object whose privileges have been affected by an `ALTER DEFAULT PRIVILEGES` command will always be shown with an explicit privilege entry that includes the effects of the `ALTER`.)

Notice that the owner's implicit grant options are not marked in the access privileges display. A `*` will appear only when grant options have been explicitly granted to someone.

The “Access privileges” column shows `(none)` when the object's privileges entry is non-null but empty. This means that no privileges are granted at all, even to the object's owner — a rare situation. (The owner still has implicit grant options in this case, and so could re-grant her own privileges; but she has none at the moment.)

5.9. Row Security Policies

Note that currently it does not work for the distributed systems. In addition to the SQL-standard [privilege system](#) available through `GRANT`, tables can have *row security policies* that restrict, on a per-user basis, which rows can be returned by normal queries or inserted, updated, or deleted by data modification commands. This feature is also known as *Row-Level Security*. By default, tables do not have any policies,

so that if a user has access privileges to a table according to the SQL privilege system, all rows within it are equally available for querying or updating.

When row security is enabled on a table (with `ALTER TABLE ... ENABLE ROW LEVEL SECURITY`), all normal access to the table for selecting rows or modifying rows must be allowed by a row security policy. (However, the table's owner is typically not subject to row security policies.) If no policy exists for the table, a default-deny policy is used, meaning that no rows are visible or can be modified. Operations that apply to the whole table, such as `TRUNCATE` and `REFERENCES`, are not subject to row security.

Row security policies can be specific to commands, or to roles, or to both. A policy can be specified to apply to `ALL` commands, or to `SELECT`, `INSERT`, `UPDATE`, or `DELETE`. Multiple roles can be assigned to a given policy, and normal role membership and inheritance rules apply.

To specify which rows are visible or modifiable according to a policy, an expression is required that returns a Boolean result. This expression will be evaluated for each row prior to any conditions or functions coming from the user's query. (The only exceptions to this rule are `leakproof` functions, which are guaranteed to not leak information; the optimizer may choose to apply such functions ahead of the row-security check.) Rows for which the expression does not return `true` will not be processed. Separate expressions may be specified to provide independent control over the rows which are visible and the rows which are allowed to be modified. Policy expressions are run as part of the query and with the privileges of the user running the query, although security-definer functions can be used to access data not available to the calling user.

Superusers and roles with the `BYPASSRLS` attribute always bypass the row security system when accessing a table. Table owners normally bypass row security as well, though a table owner can choose to be subject to row security with `ALTER TABLE ... FORCE ROW LEVEL SECURITY`.

Enabling and disabling row security, as well as adding policies to a table, is always the privilege of the table owner only.

Policies are created using the `CREATE POLICY` command, altered using the `ALTER POLICY` command, and dropped using the `DROP POLICY` command. To enable and disable row security for a given table, use the `ALTER TABLE` command.

Each policy has a name and multiple policies can be defined for a table. As policies are table-specific, each policy for a table must have a unique name. Different tables may have policies with the same name.

When multiple policies apply to a given query, they are combined using either `OR` (for permissive policies, which are the default) or using `AND` (for restrictive policies). This is similar to the rule that a given role has the privileges of all roles that they are a member of. Permissive vs. restrictive policies are discussed further below.

As a simple example, here is how to create a policy on the `account` relation to allow only members of the `managers` role to access rows, and only rows of their accounts:

```
CREATE TABLE accounts (manager text, company text, contact_email text);

ALTER TABLE accounts ENABLE ROW LEVEL SECURITY;

CREATE POLICY account_managers ON accounts TO managers
    USING (manager = current_user);
```

The policy above implicitly provides a `WITH CHECK` clause identical to its `USING` clause, so that the constraint applies both to rows selected by a command (so a manager cannot `SELECT`, `UPDATE`, or `DELETE` existing rows belonging to a different manager) and to rows modified by a command (so rows belonging to a different manager cannot be created via `INSERT` or `UPDATE`).

If no role is specified, or the special user name `PUBLIC` is used, then the policy applies to all users on the system. To allow all users to access only their own row in a `users` table, a simple policy can be used:

```
CREATE POLICY user_policy ON users
```

```
USING (user_name = current_user);
```

This works similarly to the previous example.

To use a different policy for rows that are being added to the table compared to those rows that are visible, multiple policies can be combined. This pair of policies would allow all users to view all rows in the `users` table, but only modify their own:

```
CREATE POLICY user_sel_policy ON users
  FOR SELECT
  USING (true);
CREATE POLICY user_mod_policy ON users
  USING (user_name = current_user);
```

In a `SELECT` command, these two policies are combined using `OR`, with the net effect being that all rows can be selected. In other command types, only the second policy applies, so that the effects are the same as before.

Row security can also be disabled with the `ALTER TABLE` command. Disabling row security does not remove any policies that are defined on the table; they are simply ignored. Then all rows in the table are visible and modifiable, subject to the standard SQL privileges system.

Below is a larger example of how this feature can be used in production environments. The table `passwd` emulates a Unix password file:

```
-- Simple passwd-file based example
CREATE TABLE passwd (
  user_name      text UNIQUE NOT NULL,
  pwhash         text,
  uid            int  PRIMARY KEY,
  gid            int  NOT NULL,
  real_name      text NOT NULL,
  home_phone     text,
  extra_info     text,
  home_dir       text NOT NULL,
  shell          text NOT NULL
);

CREATE ROLE admin;  -- Administrator
CREATE ROLE bob;    -- Normal user
CREATE ROLE alice;  -- Normal user

-- Populate the table
INSERT INTO passwd VALUES
  ('admin', 'xxx', 0, 0, 'Admin', '111-222-3333', null, '/root', '/bin/dash');
INSERT INTO passwd VALUES
  ('bob', 'xxx', 1, 1, 'Bob', '123-456-7890', null, '/home/bob', '/bin/zsh');
INSERT INTO passwd VALUES
  ('alice', 'xxx', 2, 1, 'Alice', '098-765-4321', null, '/home/alice', '/bin/zsh');

-- Be sure to enable row-level security on the table
ALTER TABLE passwd ENABLE ROW LEVEL SECURITY;

-- Create policies
-- Administrator can see all rows and add any rows
CREATE POLICY admin_all ON passwd TO admin USING (true) WITH CHECK (true);
-- Normal users can view all rows
CREATE POLICY all_view ON passwd FOR SELECT USING (true);
-- Normal users can update their own records, but
-- limit which shells a normal user is allowed to set
```

```
CREATE POLICY user_mod ON passwd FOR UPDATE
  USING (current_user = user_name)
  WITH CHECK (
    current_user = user_name AND
    shell IN ('/bin/bash','/bin/sh','/bin/dash','/bin/zsh','/bin/tcsh')
  );

-- Allow admin all normal rights
GRANT SELECT, INSERT, UPDATE, DELETE ON passwd TO admin;
-- Users only get select access on public columns
GRANT SELECT
  (user_name, uid, gid, real_name, home_phone, extra_info, home_dir, shell)
  ON passwd TO public;
-- Allow users to update certain columns
GRANT UPDATE
  (pwhash, real_name, home_phone, extra_info, shell)
  ON passwd TO public;
```

As with any security settings, it's important to test and ensure that the system is behaving as expected. Using the example above, this demonstrates that the permission system is working properly.

```
-- admin can view all rows and fields
postgres=> set role admin;
SET
postgres=> table passwd;
 user_name | pwhash | uid | gid | real_name | home_phone | extra_info | home_dir | shell
-----+-----+-----+-----+-----+-----+-----+-----+-----
+-----+
admin      | xxx    | 0   | 0   | Admin     | 111-222-3333 |           | /root    | /bin/dash
bob        | xxx    | 1   | 1   | Bob       | 123-456-7890 |           | /home/bob | /bin/zsh
alice      | xxx    | 2   | 1   | Alice     | 098-765-4321 |           | /home/alice | /bin/zsh
(3 rows)

-- Test what Alice is able to do
postgres=> set role alice;
SET
postgres=> table passwd;
ERROR:  permission denied for table passwd
postgres=> select user_name,real_name,home_phone,extra_info,home_dir,shell from passwd;
 user_name | real_name | home_phone | extra_info | home_dir | shell
-----+-----+-----+-----+-----+-----
admin      | Admin     | 111-222-3333 |           | /root    | /bin/dash
bob        | Bob       | 123-456-7890 |           | /home/bob | /bin/zsh
alice      | Alice     | 098-765-4321 |           | /home/alice | /bin/zsh
(3 rows)

postgres=> update passwd set user_name = 'joe';
ERROR:  permission denied for table passwd
-- Alice is allowed to change her own real_name, but no others
postgres=> update passwd set real_name = 'Alice Doe';
UPDATE 1
postgres=> update passwd set real_name = 'John Doe' where user_name = 'admin';
UPDATE 0
postgres=> update passwd set shell = '/bin/xx';
ERROR:  new row violates WITH CHECK OPTION for "passwd"
```

```
postgres=> delete from passwd;
ERROR:  permission denied for table passwd
postgres=> insert into passwd (user_name) values ('xxx');
ERROR:  permission denied for table passwd
-- Alice can change her own password; RLS silently prevents updating other rows
postgres=> update passwd set pwhash = 'abc';
UPDATE 1
```

All of the policies constructed thus far have been permissive policies, meaning that when multiple policies are applied they are combined using the “OR” Boolean operator. While permissive policies can be constructed to only allow access to rows in the intended cases, it can be simpler to combine permissive policies with restrictive policies (which the records must pass and which are combined using the “AND” Boolean operator). Building on the example above, we add a restrictive policy to require the administrator to be connected over a local Unix socket to access the records of the `passwd` table:

```
CREATE POLICY admin_local_only ON passwd AS RESTRICTIVE TO admin
    USING (pg_catalog.inet_client_addr() IS NULL);
```

We can then see that an administrator connecting over a network will not see any records, due to the restrictive policy:

```
=> SELECT current_user;
    current_user
-----
    admin
(1 row)

=> select inet_client_addr();
    inet_client_addr
-----
    127.0.0.1
(1 row)

=> TABLE passwd;
 user_name | pwhash | uid | gid | real_name | home_phone | extra_info | home_dir |
    shell
-----+-----+-----+-----+-----+-----+-----+-----+-----
(0 rows)

=> UPDATE passwd set pwhash = NULL;
UPDATE 0
```

Referential integrity checks, such as unique or primary key constraints and foreign key references, always bypass row security to ensure that data integrity is maintained. Care must be taken when developing schemas and row level policies to avoid “covert channel” leaks of information through such referential integrity checks.

In some contexts it is important to be sure that row security is not being applied. For example, when taking a backup, it could be disastrous if row security silently caused some rows to be omitted from the backup. In such a situation, you can set the [row_security](#) configuration parameter to `off`. This does not in itself bypass row security; what it does is throw an error if any query's results would get filtered by a policy. The reason for the error can then be investigated and fixed.

In the examples above, the policy expressions consider only the current values in the row to be accessed or updated. This is the simplest and best-performing case; when possible, it's best to design row security applications to work this way. If it is necessary to consult other rows or other tables to make a policy decision, that can be accomplished using sub-SELECTs, or functions that contain SELECTs, in the policy expressions. Be aware however that such accesses can create race conditions that could allow information leakage if care is not taken. As an example, consider the following table design:

```
-- definition of privilege groups
CREATE TABLE groups (group_id int PRIMARY KEY,
                     group_name text NOT NULL);

INSERT INTO groups VALUES
  (1, 'low'),
  (2, 'medium'),
  (5, 'high');

GRANT ALL ON groups TO alice; -- alice is the administrator
GRANT SELECT ON groups TO public;

-- definition of users' privilege levels
CREATE TABLE users (user_name text PRIMARY KEY,
                    group_id int NOT NULL REFERENCES groups);

INSERT INTO users VALUES
  ('alice', 5),
  ('bob', 2),
  ('mallory', 2);

GRANT ALL ON users TO alice;
GRANT SELECT ON users TO public;

-- table holding the information to be protected
CREATE TABLE information (info text,
                          group_id int NOT NULL REFERENCES groups);

INSERT INTO information VALUES
  ('barely secret', 1),
  ('slightly secret', 2),
  ('very secret', 5);

ALTER TABLE information ENABLE ROW LEVEL SECURITY;

-- a row should be visible to/updatable by users whose security group_id is
-- greater than or equal to the row's group_id
CREATE POLICY fp_s ON information FOR SELECT
  USING (group_id <= (SELECT group_id FROM users WHERE user_name = current_user));
CREATE POLICY fp_u ON information FOR UPDATE
  USING (group_id <= (SELECT group_id FROM users WHERE user_name = current_user));

-- we rely only on RLS to protect the information table
GRANT ALL ON information TO public;
```

Now suppose that `alice` wishes to change the “slightly secret” information, but decides that `mallory` should not be trusted with the new content of that row, so she does:

```
BEGIN;
UPDATE users SET group_id = 1 WHERE user_name = 'mallory';
UPDATE information SET info = 'secret from mallory' WHERE group_id = 2;
COMMIT;
```

That looks safe; there is no window wherein `mallory` should be able to see the “secret from mallory” string. However, there is a race condition here. If `mallory` is concurrently doing, say,

```
SELECT * FROM information WHERE group_id = 2 FOR UPDATE;
```

and her transaction is in `READ COMMITTED` mode, it is possible for her to see “secret from mallory”. That happens if her transaction reaches the `information` row just after `alice`'s does. It blocks waiting for

`alice`'s transaction to commit, then fetches the updated row contents thanks to the `FOR UPDATE` clause. However, it does *not* fetch an updated row for the implicit `SELECT` from `users`, because that sub-`SELECT` did not have `FOR UPDATE`; instead the `users` row is read with the snapshot taken at the start of the query. Therefore, the policy expression tests the old value of `mallory`'s privilege level and allows her to see the updated row.

There are several ways around this problem. One simple answer is to use `SELECT ... FOR SHARE` in sub-`SELECT`s in row security policies. However, that requires granting `UPDATE` privilege on the referenced table (here `users`) to the affected users, which might be undesirable. (But another row security policy could be applied to prevent them from actually exercising that privilege; or the sub-`SELECT` could be embedded into a security definer function.) Also, heavy concurrent use of row share locks on the referenced table could pose a performance problem, especially if updates of it are frequent. Another solution, practical if updates of the referenced table are infrequent, is to take an `ACCESS EXCLUSIVE` lock on the referenced table when updating it, so that no concurrent transactions could be examining old row values. Or one could just wait for all concurrent transactions to end after committing an update of the referenced table and before making changes that rely on the new security situation.

For additional details see [CREATE POLICY](#) and [ALTER TABLE](#).

5.10. Non-Global Schemas

A Postgres Pro database cluster contains one or more named databases. Roles and a few other object types are shared across the entire cluster. A client connection to the server can only access data in a single database, the one specified in the connection request.

Note

Users of a cluster do not necessarily have the privilege to access every database in the cluster. Sharing of role names means that there cannot be different roles named, say, `joe` in two databases in the same cluster; but the system can be configured to allow `joe` access to only some of the databases.

A database contains one or more named *schemas*, which in turn contain tables. Schemas also contain other kinds of named objects, including data types, functions, and operators. Within one schema, two objects of the same type cannot have the same name. Furthermore, tables, sequences, indexes, views, materialized views, and foreign tables share the same namespace, so that, for example, an index and a table must have different names if they are in the same schema. The same object name can be used in different schemas without conflict; for example, both `schema1` and `myschema` can contain tables named `mytable`. Unlike databases, schemas are not rigidly separated: a user can access objects in any of the schemas in the database they are connected to, if they have privileges to do so.

There are several reasons why one might want to use schemas:

- To allow many users to use one database without interfering with each other.
- To organize database objects into logical groups to make them more manageable.
- Third-party applications can be put into separate schemas so they do not collide with the names of other objects.

Schemas are analogous to directories at the operating system level, except that schemas cannot be nested.

First of all, when transitioning from a regular database schema to the distributed one, it makes sense to start the design with deciding how the data will be distributed in the Postgres Pro Shardman cluster. Postgres Pro Shardman distributes table rows across shards according to the hash value of the column to use for the table partitioning. In other words, the desired distribution must be even, and it aims to distribute equal parts of the data across cluster nodes and evenly distribute the workload.

When a database architect chooses the column to use for the table partitioning, the majority of typical queries executed must be taken into account to ensure the maximum performance.

In general, for most queries, especially, for those that use joins, the sharding key must be included in the query text. Otherwise, Postgres Pro Shardman will not push down queries to cluster nodes for execution, which will cause essential performance degradation as compared to usage of a single instance.

Secondly, when choosing a sharding key, it is important that it does not change. A resharding operation, that is, a change of the sharding key, is pretty time-consuming and resource-intensive. At present, Postgres Pro Shardman lacks techniques that automate this procedure. In general, if resharding is required, the data in all the sharded tables should be either moved to local tables or to sharded tables with another sharding key. Then you will have to create new sharded tables with a new sharding key and move the data back. This operation is very expensive and resource-intensive. Such operations often cannot be performed without the system outage during the migration.

Another point is that distributed transactions, that is, those that update data on several cluster shards at the same point in time, cannot be performed for free. So the better data is located and computations are performed inside one shard, the faster queries are executed. In general, the proportion of distributed and non-distributed transactions must be shifted towards non-distributed ones. Only apply distributed transactions if you have a compelling need to do it.

A Postgres Pro Shardman cluster uses two main ways to store data: sharded tables is the main way, designed for big data, and global tables, designed for small dictionaries. A sharded table contains different parts of the data in each shard, while a global table contains the same data in all shards. Efficient query execution on a Postgres Pro Shardman cluster requires that the data is properly distributed across cluster shards and primarily, a sharding key is properly selected.

5.10.1. Creating a Schema

To create a schema, use the [CREATE SCHEMA](#) command. Give the schema a name of your choice. For example:

```
CREATE SCHEMA myschema;
```

To create or access objects in a schema, write a *qualified name* consisting of the schema name and table name separated by a dot:

```
schema.table
```

This works anywhere a table name is expected, including the table modification commands and the data access commands discussed in the following chapters. (For brevity we will speak of tables only, but the same ideas apply to other kinds of named objects, such as types and functions.)

Actually, the even more general syntax

```
database.schema.table
```

can be used too, but at present this is just for pro forma compliance with the SQL standard. If you write a database name, it must be the same as the database you are connected to.

So to create a table in the new schema, use:

```
CREATE TABLE myschema.mytable (  
    ...  
);
```

To drop a schema if it's empty (all objects in it have been dropped), use:

```
DROP SCHEMA myschema;
```

To drop a schema including all contained objects, use:

```
DROP SCHEMA myschema CASCADE;
```

See [Section 5.15](#) for a description of the general mechanism behind this.

Often you will want to create a schema owned by someone else (since this is one of the ways to restrict the activities of your users to well-defined namespaces). The syntax for that is:

```
CREATE SCHEMA schema_name AUTHORIZATION user_name;
```

You can even omit the schema name, in which case the schema name will be the same as the user name. See [Section 5.10.7](#) for how this can be useful.

Schema names beginning with `pg_` are reserved for system purposes and cannot be created by users.

5.10.2. The Public Schema

In the previous sections we created tables without specifying any schema names. By default such tables (and other objects) are automatically put into a schema named “public”. Every new database contains such a schema. Thus, the following are equivalent:

```
CREATE TABLE products ( ... );
```

and:

```
CREATE TABLE public.products ( ... );
```

5.10.3. The `vault` Schema

Postgres Pro provides an enhanced security mechanism that allows you to protect sensitive data against unauthorized access of malicious users who can read or even modify such data and stay undetected. This is achieved by creating the `vault` schema and designating a separate user called security officer who manages access to the schema and its objects. In this case, the schema owner is only responsible for the management of schema objects. You must be a superuser to set the security officer for a schema. To set the roles and create the `vault` schema, follow the procedure below:

1. Create a user for the `vault` schema owner and remember to set a password for this user.

```
CREATE USER vault_owner WITH LOGIN;  
GRANT CONNECT ON DATABASE database_name TO vault_owner;
```

2. Create a user for the security officer of the `vault` schema and remember to set a password for this user.

```
CREATE USER vault_security_officer WITH LOGIN;  
GRANT CONNECT ON DATABASE database_name TO vault_security_officer;
```

3. Create the `vault` schema, assign its owner, and set the security officer.

```
CREATE SCHEMA vault_name;  
ALTER SCHEMA vault_name OWNER TO vault_owner;  
ALTER SCHEMA vault_name SECURITY OFFICER TO vault_security_officer;
```

The security officer is the only non-superuser who can `GRANT` or `REVOKE` privileges to access the `vault` schema and its objects, more specifically `pg_class`, `pg_proc`, `pg_type`, and `pg_collation`. This cannot be done by the schema owner or the object owner.

Postgres Pro also enables you to create the `vault` schema out of the existing schema by setting the security officer for it. To set the security officer, use the [ALTER SCHEMA](#) command.

Once the `vault` schema, its owner, and security officer are created, the superuser login is disabled. After that, only the schema owner and the schema security officer can manage the data and access to it, respectively. Also, no other user of the server has `ADMIN OPTION` on them.

5.10.4. The Schema Search Path

Qualified names are tedious to write, and it's often best not to wire a particular schema name into applications anyway. Therefore tables are often referred to by *unqualified names*, which consist of just the table name. The system determines which table is meant by following a *search path*, which is a list of schemas to look in. The first matching table in the search path is taken to be the one wanted. If there

is no match in the search path, an error is reported, even if matching table names exist in other schemas in the database.

The ability to create like-named objects in different schemas complicates writing a query that references precisely the same objects every time. It also opens up the potential for users to change the behavior of other users' queries, maliciously or accidentally. Due to the prevalence of unqualified names in queries and their use in Postgres Pro internals, adding a schema to `search_path` effectively trusts all users having `CREATE` privilege on that schema. When you run an ordinary query, a malicious user able to create objects in a schema of your search path can take control and execute arbitrary SQL functions as though you executed them.

The first schema named in the search path is called the current schema. Aside from being the first schema searched, it is also the schema in which new tables will be created if the `CREATE TABLE` command does not specify a schema name.

To show the current search path, use the following command:

```
SHOW search_path;
```

In the default setup this returns:

```
search_path
-----
"$user", public
```

The first element specifies that a schema with the same name as the current user is to be searched. If no such schema exists, the entry is ignored. The second element refers to the public schema that we have seen already.

The first schema in the search path that exists is the default location for creating new objects. That is the reason that by default objects are created in the public schema. When objects are referenced in any other context without schema qualification (table modification, data modification, or query commands) the search path is traversed until a matching object is found. Therefore, in the default configuration, any unqualified access again can only refer to the public schema.

To put our new schema in the path, we use:

```
SET search_path TO myschema,public;
```

(We omit the `$user` here because we have no immediate need for it.) And then we can access the table without schema qualification:

```
DROP TABLE mytable;
```

Also, since `myschema` is the first element in the path, new objects would by default be created in it.

We could also have written:

```
SET search_path TO myschema;
```

Then we no longer have access to the public schema without explicit qualification. There is nothing special about the public schema except that it exists by default. It can be dropped, too.

See also [Section 9.27](#) for other ways to manipulate the schema search path.

The search path works in the same way for data type names, function names, and operator names as it does for table names. Data type and function names can be qualified in exactly the same way as table names. If you need to write a qualified operator name in an expression, there is a special provision: you must write

```
OPERATOR(schema.operator)
```

This is needed to avoid syntactic ambiguity. An example is:

```
SELECT 3 OPERATOR(pg_catalog.+) 4;
```

In practice one usually relies on the search path for operators, so as not to have to write anything so ugly as that.

5.10.5. Schemas and Privileges

By default, users cannot access any objects in schemas they do not own. To allow that, the owner of the schema must grant the `USAGE` privilege on the schema. By default, everyone has that privilege on the schema `public`. To allow users to make use of the objects in a schema, additional privileges might need to be granted, as appropriate for the object.

A user can also be allowed to create objects in someone else's schema. To allow that, the `CREATE` privilege on the schema needs to be granted. In databases upgraded from PostgreSQL 14 or earlier, everyone has that privilege on the schema `public`. Some [usage patterns](#) call for revoking that privilege:

```
REVOKE CREATE ON SCHEMA public FROM PUBLIC;
```

(The first “public” is the schema, the second “public” means “every user”. In the first sense it is an identifier, in the second sense it is a key word, hence the different capitalization; recall the guidelines from [Section 4.1.1.](#))

5.10.6. The System Catalog Schema

In addition to `public` and user-created schemas, each database contains a `pg_catalog` schema, which contains the system tables and all the built-in data types, functions, and operators. `pg_catalog` is always effectively part of the search path. If it is not named explicitly in the path then it is implicitly searched *before* searching the path's schemas. This ensures that built-in names will always be findable. However, you can explicitly place `pg_catalog` at the end of your search path if you prefer to have user-defined names override built-in names.

Since system table names begin with `pg_`, it is best to avoid such names to ensure that you won't suffer a conflict if some future version defines a system table named the same as your table. (With the default search path, an unqualified reference to your table name would then be resolved as the system table instead.) System tables will continue to follow the convention of having names beginning with `pg_`, so that they will not conflict with unqualified user-table names so long as users avoid the `pg_` prefix.

5.10.7. Usage Patterns

Schemas can be used to organize your data in many ways. A *secure schema usage pattern* prevents untrusted users from changing the behavior of other users' queries. When a database does not use a secure schema usage pattern, users wishing to securely query that database would take protective action at the beginning of each session. Specifically, they would begin each session by setting `search_path` to the empty string or otherwise removing schemas that are writable by non-superusers from `search_path`. There are a few usage patterns easily supported by the default configuration:

- Constrain ordinary users to user-private schemas. To implement this pattern, first ensure that no schemas have public `CREATE` privileges. Then, for every user needing to create non-temporary objects, create a schema with the same name as that user, for example `CREATE SCHEMA alice AUTHORIZATION alice`. (Recall that the default search path starts with `$user`, which resolves to the user name. Therefore, if each user has a separate schema, they access their own schemas by default.) This pattern is a secure schema usage pattern unless an untrusted user is the database owner or has been granted `ADMIN OPTION` on a relevant role, in which case no secure schema usage pattern exists.

In PostgreSQL 15 and later, the default configuration supports this usage pattern. In prior versions, or when using a database that has been upgraded from a prior version, you will need to remove the public `CREATE` privilege from the `public` schema (issue `REVOKE CREATE ON SCHEMA public FROM PUBLIC`). Then consider auditing the `public` schema for objects named like objects in schema `pg_catalog`.

- Remove the `public` schema from the default search path, by modifying [postgresql.conf](#) or by issuing `ALTER ROLE ALL SET search_path = '$user'`. Then, grant privileges to create in the `public`

schema. Only qualified names will choose public schema objects. While qualified table references are fine, calls to functions in the public schema [will be unsafe or unreliable](#). If you create functions or extensions in the public schema, use the first pattern instead. Otherwise, like the first pattern, this is secure unless an untrusted user is the database owner or has been granted `ADMIN OPTION` on a relevant role.

- Keep the default search path, and grant privileges to create in the public schema. All users access the public schema implicitly. This simulates the situation where schemas are not available at all, giving a smooth transition from the non-schema-aware world. However, this is never a secure pattern. It is acceptable only when the database has a single user or a few mutually-trusting users. In databases upgraded from PostgreSQL 14 or earlier, this is the default.

For any pattern, to install shared applications (tables to be used by everyone, additional functions provided by third parties, etc.), put them into separate schemas. Remember to grant appropriate privileges to allow the other users to access them. Users can then refer to these additional objects by qualifying the names with a schema name, or they can put the additional schemas into their search path, as they choose.

5.10.8. Portability

In the SQL standard, the notion of objects in the same schema being owned by different users does not exist. Moreover, some implementations do not allow you to create schemas that have a different name than their owner. In fact, the concepts of schema and user are nearly equivalent in a database system that implements only the basic schema support specified in the standard. Therefore, many users consider qualified names to really consist of `user_name.table_name`. This is how Postgres Pro will effectively behave if you create a per-user schema for every user.

Also, there is no concept of a `public` schema in the SQL standard. For maximum conformance to the standard, you should not use the `public` schema.

Of course, some SQL database systems might not implement schemas at all, or provide namespace support by allowing (possibly limited) cross-database access. If you need to work with those systems, then maximum portability would be achieved by not using schemas at all.

5.11. Inheritance

Postgres Pro implements table inheritance, which can be a useful tool for database designers. (SQL:1999 and later define a type inheritance feature, which differs in many respects from the features described here.) It has some limitations in the distributed systems, described for each command (see [CREATE TABLE](#) and [ALTER TABLE](#)).

Let's start with an example: suppose we are trying to build a data model for cities. Each state has many cities, but only one capital. We want to be able to quickly retrieve the capital city for any particular state. This can be done by creating two tables, one for state capitals and one for cities that are not capitals. However, what happens when we want to ask for data about a city, regardless of whether it is a capital or not? The inheritance feature can help to resolve this problem. We define the `capitals` table so that it inherits from `cities`:

```
CREATE TABLE cities (  
    name          text,  
    population    float,  
    elevation     int      -- in feet  
);  
  
CREATE TABLE capitals (  
    state         char(2)  
) INHERITS (cities);
```

In this case, the `capitals` table *inherits* all the columns of its parent table, `cities`. State capitals also have an extra column, `state`, that shows their state.

In Postgres Pro, a table can inherit from zero or more other tables, and a query can reference either all rows of a table or all rows of a table plus all of its descendant tables. The latter behavior is the default. For example, the following query finds the names of all cities, including state capitals, that are located at an elevation over 500 feet:

```
SELECT name, elevation
FROM cities
WHERE elevation > 500;
```

Given the sample data from the Postgres Pro tutorial (see [Section 2.1](#)), this returns:

name	elevation
Las Vegas	2174
Mariposa	1953
Madison	845

On the other hand, the following query finds all the cities that are not state capitals and are situated at an elevation over 500 feet:

```
SELECT name, elevation
FROM ONLY cities
WHERE elevation > 500;
```

name	elevation
Las Vegas	2174
Mariposa	1953

Here the **ONLY** keyword indicates that the query should apply only to `cities`, and not any tables below `cities` in the inheritance hierarchy. Many of the commands that we have already discussed — `SELECT`, `UPDATE` and `DELETE` — support the **ONLY** keyword.

You can also write the table name with a trailing `*` to explicitly specify that descendant tables are included:

```
SELECT name, elevation
FROM cities*
WHERE elevation > 500;
```

Writing `*` is not necessary, since this behavior is always the default. However, this syntax is still supported for compatibility with older releases where the default could be changed.

In some cases you might wish to know which table a particular row originated from. There is a system column called `tableoid` in each table which can tell you the originating table:

```
SELECT c.tableoid, c.name, c.elevation
FROM cities c
WHERE c.elevation > 500;
```

which returns:

tableoid	name	elevation
139793	Las Vegas	2174
139793	Mariposa	1953
139798	Madison	845

(If you try to reproduce this example, you will probably get different numeric OIDs.) By doing a join with `pg_class` you can see the actual table names:

```
SELECT p.relname, c.name, c.elevation
FROM cities c, pg_class p
WHERE c.elevation > 500 AND c.tableoid = p.oid;
```

which returns:

relname	name	elevation
cities	Las Vegas	2174
cities	Mariposa	1953
capitals	Madison	845

Another way to get the same effect is to use the `regclass` alias type, which will print the table OID symbolically:

```
SELECT c.tableoid::regclass, c.name, c.elevation
FROM cities c
WHERE c.elevation > 500;
```

Inheritance does not automatically propagate data from `INSERT` or `COPY` commands to other tables in the inheritance hierarchy. In our example, the following `INSERT` statement will fail:

```
INSERT INTO cities (name, population, elevation, state)
VALUES ('Albany', NULL, NULL, 'NY');
```

We might hope that the data would somehow be routed to the `capitals` table, but this does not happen: `INSERT` always inserts into exactly the table specified. In some cases it is possible to redirect the insertion using a rule (see [Chapter 42](#)). However that does not help for the above case because the `cities` table does not contain the column `state`, and so the command will be rejected before the rule can be applied.

All check constraints and not-null constraints on a parent table are automatically inherited by its children, unless explicitly specified otherwise with `NO INHERIT` clauses. Other types of constraints (unique, primary key, and foreign key constraints) are not inherited.

A table can inherit from more than one parent table, in which case it has the union of the columns defined by the parent tables. Any columns declared in the child table's definition are added to these. If the same column name appears in multiple parent tables, or in both a parent table and the child's definition, then these columns are “merged” so that there is only one such column in the child table. To be merged, columns must have the same data types, else an error is raised. Inheritable check constraints and not-null constraints are merged in a similar fashion. Thus, for example, a merged column will be marked not-null if any one of the column definitions it came from is marked not-null. Check constraints are merged if they have the same name, and the merge will fail if their conditions are different.

Table inheritance is typically established when the child table is created, using the `INHERITS` clause of the `CREATE TABLE` statement. Alternatively, a table which is already defined in a compatible way can have a new parent relationship added, using the `INHERIT` variant of `ALTER TABLE`. To do this the new child table must already include columns with the same names and types as the columns of the parent. It must also include check constraints with the same names and check expressions as those of the parent. Similarly an inheritance link can be removed from a child using the `NO INHERIT` variant of `ALTER TABLE`. Dynamically adding and removing inheritance links like this can be useful when the inheritance relationship is being used for table partitioning (see [Section 5.12](#)).

One convenient way to create a compatible table that will later be made a new child is to use the `LIKE` clause in `CREATE TABLE`. This creates a new table with the same columns as the source table. If there are any `CHECK` constraints defined on the source table, the `INCLUDING CONSTRAINTS` option to `LIKE` should be specified, as the new child must have constraints matching the parent to be considered compatible.

A parent table cannot be dropped while any of its children remain. Neither can columns or check constraints of child tables be dropped or altered if they are inherited from any parent tables. If you wish to remove a table and all of its descendants, one easy way is to drop the parent table with the `CASCADE` option (see [Section 5.15](#)).

`ALTER TABLE` will propagate any changes in column data definitions and check constraints down the inheritance hierarchy. Again, dropping columns that are depended on by other tables is only possible when using the `CASCADE` option. `ALTER TABLE` follows the same rules for duplicate column merging and rejection that apply during `CREATE TABLE`.

Inherited queries perform access permission checks on the parent table only. Thus, for example, granting `UPDATE` permission on the `cities` table implies permission to update rows in the `capitals` table as well, when they are accessed through `cities`. This preserves the appearance that the data is (also) in the parent table. But the `capitals` table could not be updated directly without an additional grant. In a similar way, the parent table's row security policies (see [Section 5.9](#)) are applied to rows coming from child tables during an inherited query. A child table's policies, if any, are applied only when it is the table explicitly named in the query; and in that case, any policies attached to its parent(s) are ignored.

Foreign tables (see [Section 5.13](#)) can also be part of inheritance hierarchies, either as parent or child tables, just as regular tables can be. If a foreign table is part of an inheritance hierarchy then any operations not supported by the foreign table are not supported on the whole hierarchy either.

5.11.1. Caveats

Note that not all SQL commands are able to work on inheritance hierarchies. Commands that are used for data querying, data modification, or schema modification (e.g., `SELECT`, `UPDATE`, `DELETE`, most variants of `ALTER TABLE`, but not `INSERT` or `ALTER TABLE ... RENAME`) typically default to including child tables and support the `ONLY` notation to exclude them. Commands that do database maintenance and tuning (e.g., `REINDEX`, `VACUUM`) typically only work on individual, physical tables and do not support recursing over inheritance hierarchies. The respective behavior of each individual command is documented in its reference page ([SQL Commands](#)).

A serious limitation of the inheritance feature is that indexes (including unique constraints) and foreign key constraints only apply to single tables, not to their inheritance children. This is true on both the referencing and referenced sides of a foreign key constraint. Thus, in the terms of the above example:

- If we declared `cities.name` to be `UNIQUE` or a `PRIMARY KEY`, this would not stop the `capitals` table from having rows with names duplicating rows in `cities`. And those duplicate rows would by default show up in queries from `cities`. In fact, by default `capitals` would have no unique constraint at all, and so could contain multiple rows with the same name. You could add a unique constraint to `capitals`, but this would not prevent duplication compared to `cities`.
- Similarly, if we were to specify that `cities.name` `REFERENCES` some other table, this constraint would not automatically propagate to `capitals`. In this case you could work around it by manually adding the same `REFERENCES` constraint to `capitals`.
- Specifying that another table's column `REFERENCES cities(name)` would allow the other table to contain city names, but not capital names. There is no good workaround for this case.

Some functionality not implemented for inheritance hierarchies is implemented for declarative partitioning. Considerable care is needed in deciding whether partitioning with legacy inheritance is useful for your application.

5.12. Table Partitioning

Note that sharded tables cannot be partitioned manually, only locally. This section describes why and how to implement partitioning as part of your database design.

5.12.1. Overview

Partitioning refers to splitting what is logically one large table into smaller physical pieces. Partitioning can provide several benefits:

- Query performance can be improved dramatically in certain situations, particularly when most of the heavily accessed rows of the table are in a single partition or a small number of partitions. Partitioning effectively substitutes for the upper tree levels of indexes, making it more likely that the heavily-used parts of the indexes fit in memory.
- When queries or updates access a large percentage of a single partition, performance can be improved by using a sequential scan of that partition instead of using an index, which would require random-access reads scattered across the whole table.

- Bulk loads and deletes can be accomplished by adding or removing partitions, if the usage pattern is accounted for in the partitioning design. Dropping an individual partition using `DROP TABLE`, or doing `ALTER TABLE DETACH PARTITION`, is far faster than a bulk operation. These commands also entirely avoid the `VACUUM` overhead caused by a bulk `DELETE`.
- Seldom-used data can be migrated to cheaper and slower storage media.

These benefits will normally be worthwhile only when a table would otherwise be very large. The exact point at which a table will benefit from partitioning depends on the application, although a rule of thumb is that the size of the table should exceed the physical memory of the database server.

Postgres Pro offers built-in support for the following forms of partitioning:

Range Partitioning

The table is partitioned into “ranges” defined by a key column or set of columns, with no overlap between the ranges of values assigned to different partitions. For example, one might partition by date ranges, or by ranges of identifiers for particular business objects. Each range's bounds are understood as being inclusive at the lower end and exclusive at the upper end. For example, if one partition's range is from 1 to 10, and the next one's range is from 10 to 20, then value 10 belongs to the second partition not the first.

List Partitioning

The table is partitioned by explicitly listing which key value(s) appear in each partition.

Hash Partitioning

The table is partitioned by specifying a modulus and a remainder for each partition. Each partition will hold the rows for which the hash value of the partition key divided by the specified modulus will produce the specified remainder.

If your application needs to use other forms of partitioning not listed above, alternative methods such as inheritance and `UNION ALL` views can be used instead. Such methods offer flexibility but do not have some of the performance benefits of built-in declarative partitioning.

5.12.2. Declarative Partitioning

Postgres Pro allows you to declare that a table is divided into partitions. The table that is divided is referred to as a *partitioned table*. The declaration includes the *partitioning method* as described above, plus a list of columns or expressions to be used as the *partition key*.

The partitioned table itself is a “virtual” table having no storage of its own. Instead, the storage belongs to *partitions*, which are otherwise-ordinary tables associated with the partitioned table. Each partition stores a subset of the data as defined by its *partition bounds*. All rows inserted into a partitioned table will be routed to the appropriate one of the partitions based on the values of the partition key column(s). Updating the partition key of a row will cause it to be moved into a different partition if it no longer satisfies the partition bounds of its original partition.

Partitions may themselves be defined as partitioned tables, resulting in *sub-partitioning*. Although all partitions must have the same columns as their partitioned parent, partitions may have their own indexes, constraints and default values, distinct from those of other partitions. See [CREATE TABLE](#) for more details on creating partitioned tables and partitions.

It is not possible to turn a regular table into a partitioned table or vice versa. However, it is possible to add an existing regular or partitioned table as a partition of a partitioned table, or remove a partition from a partitioned table turning it into a standalone table; this can simplify and speed up many maintenance processes. See [ALTER TABLE](#) to learn more about the `ATTACH PARTITION` and `DETACH PARTITION` sub-commands.

Partitions can also be [foreign tables](#), although considerable care is needed because it is then the user's responsibility that the contents of the foreign table satisfy the partitioning rule. There are some other restrictions as well.

5.12.2.1. Example

Suppose we are constructing a database for a large ice cream company. The company measures peak temperatures every day as well as ice cream sales in each region. Conceptually, we want a table like:

```
CREATE TABLE measurement (  
    city_id          int not null,  
    logdate          date not null,  
    peaktemp        int,  
    unitsales       int  
);
```

We know that most queries will access just the last week's, month's or quarter's data, since the main use of this table will be to prepare online reports for management. To reduce the amount of old data that needs to be stored, we decide to keep only the most recent 3 years worth of data. At the beginning of each month we will remove the oldest month's data. In this situation we can use partitioning to help us meet all of our different requirements for the measurements table.

To use declarative partitioning in this case, use the following steps:

1. Create the `measurement` table as a partitioned table by specifying the `PARTITION BY` clause, which includes the partitioning method (`RANGE` in this case) and the list of column(s) to use as the partition key.

```
CREATE TABLE measurement (  
    city_id          int not null,  
    logdate          date not null,  
    peaktemp        int,  
    unitsales       int  
) PARTITION BY RANGE (logdate);
```

2. Create partitions. Each partition's definition must specify bounds that correspond to the partitioning method and partition key of the parent. Note that specifying bounds such that the new partition's values would overlap with those in one or more existing partitions will cause an error.

Partitions thus created are in every way normal Postgres Pro tables (or, possibly, foreign tables). It is possible to specify a tablespace and storage parameters for each partition separately.

For our example, each partition should hold one month's worth of data, to match the requirement of deleting one month's data at a time. So the commands might look like:

```
CREATE TABLE measurement_y2006m02 PARTITION OF measurement  
    FOR VALUES FROM ('2006-02-01') TO ('2006-03-01');  
  
CREATE TABLE measurement_y2006m03 PARTITION OF measurement  
    FOR VALUES FROM ('2006-03-01') TO ('2006-04-01');  
  
...  
CREATE TABLE measurement_y2007m11 PARTITION OF measurement  
    FOR VALUES FROM ('2007-11-01') TO ('2007-12-01');  
  
CREATE TABLE measurement_y2007m12 PARTITION OF measurement  
    FOR VALUES FROM ('2007-12-01') TO ('2008-01-01')  
    TABLESPACE fasttablespace;  
  
CREATE TABLE measurement_y2008m01 PARTITION OF measurement  
    FOR VALUES FROM ('2008-01-01') TO ('2008-02-01')  
    WITH (parallel_workers = 4)  
    TABLESPACE fasttablespace;
```

(Recall that adjacent partitions can share a bound value, since range upper bounds are treated as exclusive bounds.)

If you wish to implement sub-partitioning, again specify the `PARTITION BY` clause in the commands used to create individual partitions, for example:

```
CREATE TABLE measurement_y2006m02 PARTITION OF measurement
    FOR VALUES FROM ('2006-02-01') TO ('2006-03-01')
    PARTITION BY RANGE (peaktemp);
```

After creating partitions of `measurement_y2006m02`, any data inserted into `measurement` that is mapped to `measurement_y2006m02` (or data that is directly inserted into `measurement_y2006m02`, which is allowed provided its partition constraint is satisfied) will be further redirected to one of its partitions based on the `peaktemp` column. The partition key specified may overlap with the parent's partition key, although care should be taken when specifying the bounds of a sub-partition such that the set of data it accepts constitutes a subset of what the partition's own bounds allow; the system does not try to check whether that's really the case.

Inserting data into the parent table that does not map to one of the existing partitions will cause an error; an appropriate partition must be added manually.

It is not necessary to manually create table constraints describing the partition boundary conditions for partitions. Such constraints will be created automatically.

3. Create an index on the key column(s), as well as any other indexes you might want, on the partitioned table. (The key index is not strictly necessary, but in most scenarios it is helpful.) This automatically creates a matching index on each partition, and any partitions you create or attach later will also have such an index. An index or unique constraint declared on a partitioned table is “virtual” in the same way that the partitioned table is: the actual data is in child indexes on the individual partition tables.

```
CREATE INDEX ON measurement (logdate);
```

4. Ensure that the [enable_partition_pruning](#) configuration parameter is not disabled in `postgresql.conf`. If it is, queries will not be optimized as desired.

In the above example we would be creating a new partition each month, so it might be wise to write a script that generates the required DDL automatically.

5.12.2.2. Partition Maintenance

Normally the set of partitions established when initially defining the table is not intended to remain static. It is common to want to remove partitions holding old data and periodically add new partitions for new data. One of the most important advantages of partitioning is precisely that it allows this otherwise painful task to be executed nearly instantaneously by manipulating the partition structure, rather than physically moving large amounts of data around.

The simplest option for removing old data is to drop the partition that is no longer necessary:

```
DROP TABLE measurement_y2006m02;
```

This can very quickly delete millions of records because it doesn't have to individually delete every record. Note however that the above command requires taking an `ACCESS EXCLUSIVE` lock on the parent table.

Another option that is often preferable is to remove the partition from the partitioned table but retain access to it as a table in its own right. This has two forms:

```
ALTER TABLE measurement DETACH PARTITION measurement_y2006m02;
ALTER TABLE measurement DETACH PARTITION measurement_y2006m02 CONCURRENTLY;
```

These allow further operations to be performed on the data before it is dropped. For example, this is often a useful time to back up the data using `COPY`, `pg_dump`, or similar tools. It might also be a useful time to aggregate data into smaller formats, perform other data manipulations, or run reports. The first form of the command requires an `ACCESS EXCLUSIVE` lock on the parent table. Adding the `CONCURRENTLY` qualifier as in the second form allows the detach operation to require only `SHARE UPDATE EXCLUSIVE` lock on the parent table, but see [ALTER TABLE ... DETACH PARTITION](#) for details on the restrictions.

Similarly we can add a new partition to handle new data. We can create an empty partition in the partitioned table just as the original partitions were created above:

```
CREATE TABLE measurement_y2008m02 PARTITION OF measurement
    FOR VALUES FROM ('2008-02-01') TO ('2008-03-01')
    TABLESPACE fasttablespace;
```

As an alternative to creating a new partition, it is sometimes more convenient to create a new table separate from the partition structure and attach it as a partition later. This allows new data to be loaded, checked, and transformed prior to it appearing in the partitioned table. Moreover, the `ATTACH PARTITION` operation requires only a `SHARE UPDATE EXCLUSIVE` lock on the partitioned table rather than the `ACCESS EXCLUSIVE` lock required by `CREATE TABLE ... PARTITION OF`, so it is more friendly to concurrent operations on the partitioned table; see [ALTER TABLE ... ATTACH PARTITION](#) for additional details. The `CREATE TABLE ... LIKE` option can be helpful to avoid tediously repeating the parent table's definition; for example:

```
CREATE TABLE measurement_y2008m02
    (LIKE measurement INCLUDING DEFAULTS INCLUDING CONSTRAINTS)
    TABLESPACE fasttablespace;

ALTER TABLE measurement_y2008m02 ADD CONSTRAINT y2008m02
    CHECK ( logdate >= DATE '2008-02-01' AND logdate < DATE '2008-03-01' );

\copy measurement_y2008m02 from 'measurement_y2008m02'
-- possibly some other data preparation work

ALTER TABLE measurement ATTACH PARTITION measurement_y2008m02
    FOR VALUES FROM ('2008-02-01') TO ('2008-03-01' );
```

Note that when running the `ATTACH PARTITION` command, the table will be scanned to validate the partition constraint while holding an `ACCESS EXCLUSIVE` lock on that partition. As shown above, it is recommended to avoid this scan by creating a `CHECK` constraint matching the expected partition constraint on the table prior to attaching it. Once the `ATTACH PARTITION` is complete, it is recommended to drop the now-redundant `CHECK` constraint. If the table being attached is itself a partitioned table, then each of its sub-partitions will be recursively locked and scanned until either a suitable `CHECK` constraint is encountered or the leaf partitions are reached.

Similarly, if the partitioned table has a `DEFAULT` partition, it is recommended to create a `CHECK` constraint which excludes the to-be-attached partition's constraint. If this is not done, the `DEFAULT` partition will be scanned to verify that it contains no records which should be located in the partition being attached. This operation will be performed whilst holding an `ACCESS EXCLUSIVE` lock on the `DEFAULT` partition. If the `DEFAULT` partition is itself a partitioned table, then each of its partitions will be recursively checked in the same way as the table being attached, as mentioned above.

As mentioned earlier, it is possible to create indexes on partitioned tables so that they are applied automatically to the entire hierarchy. This can be very convenient as not only will all existing partitions be indexed, but any future partitions will be as well. To avoid long lock times, you can use `CREATE INDEX ON ONLY` the partitioned table, which creates the new index marked as invalid, preventing automatic application to existing partitions. Instead, indexes can then be created individually on each partition using `CONCURRENTLY` and *attached* to the partitioned index on the parent using `ALTER INDEX ... ATTACH PARTITION`. Once indexes for all the partitions are attached to the parent index, the parent index will be marked valid automatically. Example:

```
CREATE INDEX measurement_usls_idx ON ONLY measurement (unitsales);

CREATE INDEX CONCURRENTLY measurement_usls_200602_idx
    ON measurement_y2006m02 (unitsales);
ALTER INDEX measurement_usls_idx
    ATTACH PARTITION measurement_usls_200602_idx;
...
```

This technique can be used with `UNIQUE` and `PRIMARY KEY` constraints too; the indexes are created implicitly when the constraint is created. Example:

```
ALTER TABLE ONLY measurement ADD UNIQUE (city_id, logdate);

ALTER TABLE measurement_y2006m02 ADD UNIQUE (city_id, logdate);
ALTER INDEX measurement_city_id_logdate_key
  ATTACH PARTITION measurement_y2006m02_city_id_logdate_key;
...
```

There is also an option for merging multiple table partitions into a single partition using the `ALTER TABLE ... MERGE PARTITIONS`. This feature simplifies the management of partitioned tables by allowing users to combine partitions that are no longer needed as separate entities. It's important to note that this operation is not supported for hash-partitioned tables and acquires an `ACCESS EXCLUSIVE` lock, which could impact high-load systems due to the lock's restrictive nature. For example, we can merge three monthly partitions into one quarter partition:

```
ALTER TABLE measurement
  MERGE PARTITIONS (measurement_y2006m01,
                  measurement_y2006m02,
                  measurement_y2006m03) INTO measurement_y2006q1;
```

Similarly to merging multiple table partitions, there is an option for splitting a single partition into multiple using the `ALTER TABLE ... SPLIT PARTITION`. This feature could come in handy when one partition grows too big and needs to be split into multiple. It's important to note that this operation is not supported for hash-partitioned tables and acquires an `ACCESS EXCLUSIVE` lock, which could impact high-load systems due to the lock's restrictive nature. For example, we can split the quarter partition back to monthly partitions:

```
ALTER TABLE measurement SPLIT PARTITION measurement_y2006q1 INTO
  (PARTITION measurement_y2006m01 FOR VALUES FROM ('2006-01-01') TO ('2006-02-01'),
   PARTITION measurement_y2006m02 FOR VALUES FROM ('2006-02-01') TO ('2006-03-01'),
   PARTITION measurement_y2006m03 FOR VALUES FROM ('2006-03-01') TO ('2006-04-01'));
```

5.12.2.3. Limitations

The following limitations apply to partitioned tables:

- To create a unique or primary key constraint on a partitioned table, the partition keys must not include any expressions or function calls and the constraint's columns must include all of the partition key columns. This limitation exists because the individual indexes making up the constraint can only directly enforce uniqueness within their own partitions; therefore, the partition structure itself must guarantee that there are not duplicates in different partitions.
- Similarly an exclusion constraint must include all the partition key columns. Furthermore the constraint must compare those columns for equality (not e.g. `&&`). Again, this limitation stems from not being able to enforce cross-partition restrictions. The constraint may include additional columns that aren't part of the partition key, and it may compare those with any operators you like.
- `BEFORE ROW` triggers on `INSERT` cannot change which partition is the final destination for a new row.
- Mixing temporary and permanent relations in the same partition tree is not allowed. Hence, if the partitioned table is permanent, so must be its partitions and likewise if the partitioned table is temporary. When using temporary relations, all members of the partition tree have to be from the same session.

Individual partitions are linked to their partitioned table using inheritance behind-the-scenes. However, it is not possible to use all of the generic features of inheritance with declaratively partitioned tables or their partitions, as discussed below. Notably, a partition cannot have any parents other than the partitioned table it is a partition of, nor can a table inherit from both a partitioned table and a regular table. That means partitioned tables and their partitions never share an inheritance hierarchy with regular tables.

Since a partition hierarchy consisting of the partitioned table and its partitions is still an inheritance hierarchy, `tableoid` and all the normal rules of inheritance apply as described in [Section 5.11](#), with a few exceptions:

- Partitions cannot have columns that are not present in the parent. It is not possible to specify columns when creating partitions with `CREATE TABLE`, nor is it possible to add columns to partitions after-the-fact using `ALTER TABLE`. Tables may be added as a partition with `ALTER TABLE ... ATTACH PARTITION` only if their columns exactly match the parent.
- Both `CHECK` and `NOT NULL` constraints of a partitioned table are always inherited by all its partitions. `CHECK` constraints that are marked `NO INHERIT` are not allowed to be created on partitioned tables. You cannot drop a `NOT NULL` constraint on a partition's column if the same constraint is present in the parent table.
- Using `ONLY` to add or drop a constraint on only the partitioned table is supported as long as there are no partitions. Once partitions exist, using `ONLY` will result in an error for any constraints other than `UNIQUE` and `PRIMARY KEY`. Instead, constraints on the partitions themselves can be added and (if they are not present in the parent table) dropped.
- As a partitioned table does not have any data itself, attempts to use `TRUNCATE ONLY` on a partitioned table will always return an error.

5.12.3. Partitioning Using Inheritance

While the built-in declarative partitioning is suitable for most common use cases, there are some circumstances where a more flexible approach may be useful. Partitioning can be implemented using table inheritance, which allows for several features not supported by declarative partitioning, such as:

- For declarative partitioning, partitions must have exactly the same set of columns as the partitioned table, whereas with table inheritance, child tables may have extra columns not present in the parent.
- Table inheritance allows for multiple inheritance.
- Declarative partitioning only supports range, list and hash partitioning, whereas table inheritance allows data to be divided in a manner of the user's choosing. (Note, however, that if constraint exclusion is unable to prune child tables effectively, query performance might be poor.)

5.12.3.1. Example

This example builds a partitioning structure equivalent to the declarative partitioning example above. Use the following steps:

1. Create the “root” table, from which all of the “child” tables will inherit. This table will contain no data. Do not define any check constraints on this table, unless you intend them to be applied equally to all child tables. There is no point in defining any indexes or unique constraints on it, either. For our example, the root table is the `measurement` table as originally defined:

```
CREATE TABLE measurement (  
    city_id          int not null,  
    logdate          date not null,  
    peaktemp         int,  
    unitsales        int  
);
```

2. Create several “child” tables that each inherit from the root table. Normally, these tables will not add any columns to the set inherited from the root. Just as with declarative partitioning, these tables are in every way normal Postgres Pro tables (or foreign tables).

```
CREATE TABLE measurement_y2006m02 () INHERITS (measurement);  
CREATE TABLE measurement_y2006m03 () INHERITS (measurement);  
...  
CREATE TABLE measurement_y2007m11 () INHERITS (measurement);  
CREATE TABLE measurement_y2007m12 () INHERITS (measurement);
```

```
CREATE TABLE measurement_y2008m01 () INHERITS (measurement);
```

3. Add non-overlapping table constraints to the child tables to define the allowed key values in each.

Typical examples would be:

```
CHECK ( x = 1 )
CHECK ( county IN ( 'Oxfordshire', 'Buckinghamshire', 'Warwickshire' ))
CHECK ( outletID >= 100 AND outletID < 200 )
```

Ensure that the constraints guarantee that there is no overlap between the key values permitted in different child tables. A common mistake is to set up range constraints like:

```
CHECK ( outletID BETWEEN 100 AND 200 )
CHECK ( outletID BETWEEN 200 AND 300 )
```

This is wrong since it is not clear which child table the key value 200 belongs in. Instead, ranges should be defined in this style:

```
CREATE TABLE measurement_y2006m02 (
    CHECK ( logdate >= DATE '2006-02-01' AND logdate < DATE '2006-03-01' )
) INHERITS (measurement);
```

```
CREATE TABLE measurement_y2006m03 (
    CHECK ( logdate >= DATE '2006-03-01' AND logdate < DATE '2006-04-01' )
) INHERITS (measurement);
```

...

```
CREATE TABLE measurement_y2007m11 (
    CHECK ( logdate >= DATE '2007-11-01' AND logdate < DATE '2007-12-01' )
) INHERITS (measurement);
```

```
CREATE TABLE measurement_y2007m12 (
    CHECK ( logdate >= DATE '2007-12-01' AND logdate < DATE '2008-01-01' )
) INHERITS (measurement);
```

```
CREATE TABLE measurement_y2008m01 (
    CHECK ( logdate >= DATE '2008-01-01' AND logdate < DATE '2008-02-01' )
) INHERITS (measurement);
```

4. For each child table, create an index on the key column(s), as well as any other indexes you might want.

```
CREATE INDEX measurement_y2006m02_logdate ON measurement_y2006m02 (logdate);
CREATE INDEX measurement_y2006m03_logdate ON measurement_y2006m03 (logdate);
CREATE INDEX measurement_y2007m11_logdate ON measurement_y2007m11 (logdate);
CREATE INDEX measurement_y2007m12_logdate ON measurement_y2007m12 (logdate);
CREATE INDEX measurement_y2008m01_logdate ON measurement_y2008m01 (logdate);
```

5. We want our application to be able to say `INSERT INTO measurement ...` and have the data be redirected into the appropriate child table. We can arrange that by attaching a suitable trigger function to the root table. If data will be added only to the latest child, we can use a very simple trigger function:

```
CREATE OR REPLACE FUNCTION measurement_insert_trigger()
RETURNS TRIGGER AS $$
BEGIN
    INSERT INTO measurement_y2008m01 VALUES (NEW.*);
    RETURN NULL;
END;
$$
LANGUAGE plpgsql;
```

After creating the function, we create a trigger which calls the trigger function:

```
CREATE TRIGGER insert_measurement_trigger
```



```
BEFORE INSERT ON measurement
FOR EACH ROW EXECUTE FUNCTION measurement_insert_trigger();
```

We must redefine the trigger function each month so that it always inserts into the current child table. The trigger definition does not need to be updated, however.

We might want to insert data and have the server automatically locate the child table into which the row should be added. We could do this with a more complex trigger function, for example:

```
CREATE OR REPLACE FUNCTION measurement_insert_trigger()
RETURNS TRIGGER AS $$
BEGIN
    IF ( NEW.logdate >= DATE '2006-02-01' AND
        NEW.logdate < DATE '2006-03-01' ) THEN
        INSERT INTO measurement_y2006m02 VALUES (NEW.*);
    ELSIF ( NEW.logdate >= DATE '2006-03-01' AND
        NEW.logdate < DATE '2006-04-01' ) THEN
        INSERT INTO measurement_y2006m03 VALUES (NEW.*);
    ...
    ELSIF ( NEW.logdate >= DATE '2008-01-01' AND
        NEW.logdate < DATE '2008-02-01' ) THEN
        INSERT INTO measurement_y2008m01 VALUES (NEW.*);
    ELSE
        RAISE EXCEPTION 'Date out of range. Fix the measurement_insert_trigger()
function!';
    END IF;
    RETURN NULL;
END;
$$
LANGUAGE plpgsql;
```

The trigger definition is the same as before. Note that each IF test must exactly match the CHECK constraint for its child table.

While this function is more complex than the single-month case, it doesn't need to be updated as often, since branches can be added in advance of being needed.

Note

In practice, it might be best to check the newest child first, if most inserts go into that child. For simplicity, we have shown the trigger's tests in the same order as in other parts of this example.

A different approach to redirecting inserts into the appropriate child table is to set up rules, instead of a trigger, on the root table. For example:

```
CREATE RULE measurement_insert_y2006m02 AS
ON INSERT TO measurement WHERE
    ( logdate >= DATE '2006-02-01' AND logdate < DATE '2006-03-01' )
DO INSTEAD
    INSERT INTO measurement_y2006m02 VALUES (NEW.*);
...
CREATE RULE measurement_insert_y2008m01 AS
ON INSERT TO measurement WHERE
    ( logdate >= DATE '2008-01-01' AND logdate < DATE '2008-02-01' )
DO INSTEAD
    INSERT INTO measurement_y2008m01 VALUES (NEW.*);
```

A rule has significantly more overhead than a trigger, but the overhead is paid once per query rather than once per row, so this method might be advantageous for bulk-insert situations. In most cases, however, the trigger method will offer better performance.

Be aware that `COPY` ignores rules. If you want to use `COPY` to insert data, you'll need to copy into the correct child table rather than directly into the root. `COPY` does fire triggers, so you can use it normally if you use the trigger approach.

Another disadvantage of the rule approach is that there is no simple way to force an error if the set of rules doesn't cover the insertion date; the data will silently go into the root table instead.

6. Ensure that the `constraint_exclusion` configuration parameter is not disabled in `postgresql.conf`; otherwise child tables may be accessed unnecessarily.

As we can see, a complex table hierarchy could require a substantial amount of DDL. In the above example we would be creating a new child table each month, so it might be wise to write a script that generates the required DDL automatically.

5.12.3.2. Maintenance for Inheritance Partitioning

To remove old data quickly, simply drop the child table that is no longer necessary:

```
DROP TABLE measurement_y2006m02;
```

To remove the child table from the inheritance hierarchy table but retain access to it as a table in its own right:

```
ALTER TABLE measurement_y2006m02 NO INHERIT measurement;
```

To add a new child table to handle new data, create an empty child table just as the original children were created above:

```
CREATE TABLE measurement_y2008m02 (  
    CHECK ( logdate >= DATE '2008-02-01' AND logdate < DATE '2008-03-01' )  
) INHERITS (measurement);
```

Alternatively, one may want to create and populate the new child table before adding it to the table hierarchy. This could allow data to be loaded, checked, and transformed before being made visible to queries on the parent table.

```
CREATE TABLE measurement_y2008m02  
    (LIKE measurement INCLUDING DEFAULTS INCLUDING CONSTRAINTS);  
ALTER TABLE measurement_y2008m02 ADD CONSTRAINT y2008m02  
    CHECK ( logdate >= DATE '2008-02-01' AND logdate < DATE '2008-03-01' );  
\copy measurement_y2008m02 from 'measurement_y2008m02'  
-- possibly some other data preparation work  
ALTER TABLE measurement_y2008m02 INHERIT measurement;
```

5.12.3.3. Caveats

The following caveats apply to partitioning implemented using inheritance:

- There is no automatic way to verify that all of the `CHECK` constraints are mutually exclusive. It is safer to create code that generates child tables and creates and/or modifies associated objects than to write each by hand.
- Indexes and foreign key constraints apply to single tables and not to their inheritance children, hence they have some [caveats](#) to be aware of.
- The schemes shown here assume that the values of a row's key column(s) never change, or at least do not change enough to require it to move to another partition. An `UPDATE` that attempts to do that will fail because of the `CHECK` constraints. If you need to handle such cases, you can put suitable update triggers on the child tables, but it makes management of the structure much more complicated.

- If you are using manual `VACUUM` or `ANALYZE` commands, don't forget that you need to run them on each child table individually. A command like:

```
ANALYZE measurement;
```

will only process the root table.
- `INSERT` statements with `ON CONFLICT` clauses are unlikely to work as expected, as the `ON CONFLICT` action is only taken in case of unique violations on the specified target relation, not its child relations.
- Triggers or rules will be needed to route rows to the desired child table, unless the application is explicitly aware of the partitioning scheme. Triggers may be complicated to write, and will be much slower than the tuple routing performed internally by declarative partitioning.

5.12.4. Partition Pruning

Partition pruning is a query optimization technique that improves performance for declaratively partitioned tables. As an example:

```
SET enable_partition_pruning = on;           -- the default
SELECT count(*) FROM measurement WHERE logdate >= DATE '2008-01-01';
```

Without partition pruning, the above query would scan each of the partitions of the `measurement` table. With partition pruning enabled, the planner will examine the definition of each partition and prove that the partition need not be scanned because it could not contain any rows meeting the query's `WHERE` clause. When the planner can prove this, it excludes (*prunes*) the partition from the query plan.

By using the `EXPLAIN` command and the [enable_partition_pruning](#) configuration parameter, it's possible to show the difference between a plan for which partitions have been pruned and one for which they have not. A typical unoptimized plan for this type of table setup is:

```
SET enable_partition_pruning = off;
EXPLAIN SELECT count(*) FROM measurement WHERE logdate >= DATE '2008-01-01';
                                QUERY PLAN
-----
Aggregate  (cost=188.76..188.77 rows=1 width=8)
-> Append  (cost=0.00..181.05 rows=3085 width=0)
    -> Seq Scan on measurement_y2006m02 (cost=0.00..33.12 rows=617 width=0)
        Filter: (logdate >= '2008-01-01'::date)
    -> Seq Scan on measurement_y2006m03 (cost=0.00..33.12 rows=617 width=0)
        Filter: (logdate >= '2008-01-01'::date)
...
    -> Seq Scan on measurement_y2007m11 (cost=0.00..33.12 rows=617 width=0)
        Filter: (logdate >= '2008-01-01'::date)
    -> Seq Scan on measurement_y2007m12 (cost=0.00..33.12 rows=617 width=0)
        Filter: (logdate >= '2008-01-01'::date)
    -> Seq Scan on measurement_y2008m01 (cost=0.00..33.12 rows=617 width=0)
        Filter: (logdate >= '2008-01-01'::date)
```

Some or all of the partitions might use index scans instead of full-table sequential scans, but the point here is that there is no need to scan the older partitions at all to answer this query. When we enable partition pruning, we get a significantly cheaper plan that will deliver the same answer:

```
SET enable_partition_pruning = on;
EXPLAIN SELECT count(*) FROM measurement WHERE logdate >= DATE '2008-01-01';
                                QUERY PLAN
-----
Aggregate  (cost=37.75..37.76 rows=1 width=8)
-> Seq Scan on measurement_y2008m01 (cost=0.00..33.12 rows=617 width=0)
    Filter: (logdate >= '2008-01-01'::date)
```

Note that partition pruning is driven only by the constraints defined implicitly by the partition keys, not by the presence of indexes. Therefore it isn't necessary to define indexes on the key columns. Whether

an index needs to be created for a given partition depends on whether you expect that queries that scan the partition will generally scan a large part of the partition or just a small part. An index will be helpful in the latter case but not the former.

Partition pruning can be performed not only during the planning of a given query, but also during its execution. This is useful as it can allow more partitions to be pruned when clauses contain expressions whose values are not known at query planning time, for example, parameters defined in a `PREPARE` statement, using a value obtained from a subquery, or using a parameterized value on the inner side of a nested loop join. Partition pruning during execution can be performed at any of the following times:

- During initialization of the query plan. Partition pruning can be performed here for parameter values which are known during the initialization phase of execution. Partitions which are pruned during this stage will not show up in the query's `EXPLAIN` or `EXPLAIN ANALYZE`. It is possible to determine the number of partitions which were removed during this phase by observing the “Subplans Removed” property in the `EXPLAIN` output. It's important to note that any partitions removed by the partition pruning done at this stage are still locked at the beginning of execution.
- During actual execution of the query plan. Partition pruning may also be performed here to remove partitions using values which are only known during actual query execution. This includes values from subqueries and values from execution-time parameters such as those from parameterized nested loop joins. Since the value of these parameters may change many times during the execution of the query, partition pruning is performed whenever one of the execution parameters being used by partition pruning changes. Determining if partitions were pruned during this phase requires careful inspection of the `loops` property in the `EXPLAIN ANALYZE` output. Subplans corresponding to different partitions may have different values for it depending on how many times each of them was pruned during execution. Some may be shown as `(never executed)` if they were pruned every time.

Partition pruning can be disabled using the [enable_partition_pruning](#) setting.

5.12.5. Partitioning and Constraint Exclusion

Constraint exclusion is a query optimization technique similar to partition pruning. While it is primarily used for partitioning implemented using the legacy inheritance method, it can be used for other purposes, including with declarative partitioning.

Constraint exclusion works in a very similar way to partition pruning, except that it uses each table's `CHECK` constraints — which gives it its name — whereas partition pruning uses the table's partition bounds, which exist only in the case of declarative partitioning. Another difference is that constraint exclusion is only applied at plan time; there is no attempt to remove partitions at execution time.

The fact that constraint exclusion uses `CHECK` constraints, which makes it slow compared to partition pruning, can sometimes be used as an advantage: because constraints can be defined even on declaratively-partitioned tables, in addition to their internal partition bounds, constraint exclusion may be able to elide additional partitions from the query plan.

The default (and recommended) setting of [constraint_exclusion](#) is neither `on` nor `off`, but an intermediate setting called `partition`, which causes the technique to be applied only to queries that are likely to be working on inheritance partitioned tables. The `on` setting causes the planner to examine `CHECK` constraints in all queries, even simple ones that are unlikely to benefit.

The following caveats apply to constraint exclusion:

- Constraint exclusion is only applied during query planning, unlike partition pruning, which can also be applied during query execution.
- Constraint exclusion only works when the query's `WHERE` clause contains constants (or externally supplied parameters). For example, a comparison against a non-immutable function such as `CURRENT_TIMESTAMP` cannot be optimized, since the planner cannot know which child table the function's value might fall into at run time.

- Keep the partitioning constraints simple, else the planner may not be able to prove that child tables might not need to be visited. Use simple equality conditions for list partitioning, or simple range tests for range partitioning, as illustrated in the preceding examples. A good rule of thumb is that partitioning constraints should contain only comparisons of the partitioning column(s) to constants using B-tree-indexable operators, because only B-tree-indexable column(s) are allowed in the partition key.
- All constraints on all children of the parent table are examined during constraint exclusion, so large numbers of children are likely to increase query planning time considerably. So the legacy inheritance based partitioning will work well with up to perhaps a hundred child tables; don't try to use many thousands of children.

5.12.6. Best Practices for Declarative Partitioning

The choice of how to partition a table should be made carefully, as the performance of query planning and execution can be negatively affected by poor design.

One of the most critical design decisions will be the column or columns by which you partition your data. Often the best choice will be to partition by the column or set of columns which most commonly appear in `WHERE` clauses of queries being executed on the partitioned table. `WHERE` clauses that are compatible with the partition bound constraints can be used to prune unneeded partitions. However, you may be forced into making other decisions by requirements for the `PRIMARY KEY` or a `UNIQUE` constraint. Removal of unwanted data is also a factor to consider when planning your partitioning strategy. An entire partition can be detached fairly quickly, so it may be beneficial to design the partition strategy in such a way that all data to be removed at once is located in a single partition.

Choosing the target number of partitions that the table should be divided into is also a critical decision to make. Not having enough partitions may mean that indexes remain too large and that data locality remains poor which could result in low cache hit ratios. However, dividing the table into too many partitions can also cause issues. Too many partitions can mean longer query planning times and higher memory consumption during both query planning and execution, as further described below. When choosing how to partition your table, it's also important to consider what changes may occur in the future. For example, if you choose to have one partition per customer and you currently have a small number of large customers, consider the implications if in several years you instead find yourself with a large number of small customers. In this case, it may be better to choose to partition by `HASH` and choose a reasonable number of partitions rather than trying to partition by `LIST` and hoping that the number of customers does not increase beyond what it is practical to partition the data by.

Sub-partitioning can be useful to further divide partitions that are expected to become larger than other partitions. Another option is to use range partitioning with multiple columns in the partition key. Either of these can easily lead to excessive numbers of partitions, so restraint is advisable.

It is important to consider the overhead of partitioning during query planning and execution. The query planner is generally able to handle partition hierarchies with up to a few thousand partitions fairly well, provided that typical queries allow the query planner to prune all but a small number of partitions. Planning times become longer and memory consumption becomes higher when more partitions remain after the planner performs partition pruning. Another reason to be concerned about having a large number of partitions is that the server's memory consumption may grow significantly over time, especially if many sessions touch large numbers of partitions. That's because each partition requires its metadata to be loaded into the local memory of each session that touches it.

With data warehouse type workloads, it can make sense to use a larger number of partitions than with an OLTP type workload. Generally, in data warehouses, query planning time is less of a concern as the majority of processing time is spent during query execution. With either of these two types of workload, it is important to make the right decisions early, as re-partitioning large quantities of data can be painfully slow. Simulations of the intended workload are often beneficial for optimizing the partitioning strategy. Never just assume that more partitions are better than fewer partitions, nor vice-versa.

5.13. Foreign Data

Postgres Pro implements portions of the SQL/MED specification, allowing you to access data that resides outside Postgres Pro using regular SQL queries. Such data is referred to as *foreign data*. (Note that this usage is not to be confused with foreign keys, which are a type of constraint within the database.)

Foreign data is accessed with help from a *foreign data wrapper*. A foreign data wrapper is a library that can communicate with an external data source, hiding the details of connecting to the data source and obtaining data from it. There are some foreign data wrappers available as `contrib` modules; see [Appendix F](#). Other kinds of foreign data wrappers might be found as third party products. If none of the existing foreign data wrappers suit your needs, you can write your own; see [Chapter 58](#).

To access foreign data, you need to create a *foreign server* object, which defines how to connect to a particular external data source according to the set of options used by its supporting foreign data wrapper. Then you need to create one or more *foreign tables*, which define the structure of the remote data. A foreign table can be used in queries just like a normal table, but a foreign table has no storage in the Postgres Pro server. Whenever it is used, Postgres Pro asks the foreign data wrapper to fetch data from the external source, or transmit data to the external source in the case of update commands.

Accessing remote data may require authenticating to the external data source. This information can be provided by a *user mapping*, which can provide additional data such as user names and passwords based on the current Postgres Pro role.

5.14. Other Database Objects

Tables are the central objects in a relational database structure, because they hold your data. But they are not the only objects that exist in a database. Many other kinds of objects can be created to make the use and management of the data more efficient or convenient. They are not discussed in this chapter, but we give you a list here so that you are aware of what is possible:

- Views
- Functions, procedures, and operators
- Data types and domains
- Triggers and rewrite rules

Detailed information on these topics appears in [Part V](#).

5.15. Dependency Tracking

When you create complex database structures involving many tables with foreign key constraints, views, triggers, functions, etc. you implicitly create a net of dependencies between the objects. For instance, a table with a foreign key constraint depends on the table it references.

To ensure the integrity of the entire database structure, Postgres Pro makes sure that you cannot drop objects that other objects still depend on. For example, attempting to drop the `products` table we considered in [Section 5.5.5](#), with the `orders` table depending on it, would result in an error message like this:

```
DROP TABLE products;
```

```
ERROR:  cannot drop table products because other objects depend on it
DETAIL:  constraint orders_product_no_fkey on table orders depends on table products
HINT:   Use DROP ... CASCADE to drop the dependent objects too.
```

The error message contains a useful hint: if you do not want to bother deleting all the dependent objects individually, you can run:

```
DROP TABLE products CASCADE;
```

and all the dependent objects will be removed, as will any objects that depend on them, recursively. In this case, it doesn't remove the `orders` table, it only removes the foreign key constraint. It stops there because nothing depends on the foreign key constraint. (If you want to check what `DROP ... CASCADE` will do, run `DROP` without `CASCADE` and read the `DETAIL` output.)

Almost all `DROP` commands in Postgres Pro support specifying `CASCADE`. Of course, the nature of the possible dependencies varies with the type of the object. You can also write `RESTRICT` instead of `CASCADE` to get the default behavior, which is to prevent dropping objects that any other objects depend on.

Note

According to the SQL standard, specifying either `RESTRICT` or `CASCADE` is required in a `DROP` command. No database system actually enforces that rule, but whether the default behavior is `RESTRICT` or `CASCADE` varies across systems.

If a `DROP` command lists multiple objects, `CASCADE` is only required when there are dependencies outside the specified group. For example, when saying `DROP TABLE tab1, tab2` the existence of a foreign key referencing `tab1` from `tab2` would not mean that `CASCADE` is needed to succeed.

For a user-defined function or procedure whose body is defined as a string literal, Postgres Pro tracks dependencies associated with the function's externally-visible properties, such as its argument and result types, but *not* dependencies that could only be known by examining the function body. As an example, consider this situation:

```
CREATE TYPE rainbow AS ENUM ('red', 'orange', 'yellow',
                             'green', 'blue', 'purple');

CREATE TABLE my_colors (color rainbow, note text);

CREATE FUNCTION get_color_note (rainbow) RETURNS text AS
    'SELECT note FROM my_colors WHERE color = $1'
LANGUAGE SQL;
```

(See [Section 39.5](#) for an explanation of SQL-language functions.) Postgres Pro will be aware that the `get_color_note` function depends on the `rainbow` type: dropping the type would force dropping the function, because its argument type would no longer be defined. But Postgres Pro will not consider `get_color_note` to depend on the `my_colors` table, and so will not drop the function if the table is dropped. While there are disadvantages to this approach, there are also benefits. The function is still valid in some sense if the table is missing, though executing it would cause an error; creating a new table of the same name would allow the function to work again.

On the other hand, for an SQL-language function or procedure whose body is written in SQL-standard style, the body is parsed at function definition time and all dependencies recognized by the parser are stored. Thus, if we write the function above as

```
CREATE FUNCTION get_color_note (rainbow) RETURNS text
BEGIN ATOMIC
    SELECT note FROM my_colors WHERE color = $1;
END;
```

then the function's dependency on the `my_colors` table will be known and enforced by `DROP`.

Chapter 6. Data Manipulation

The previous chapter discussed how to create tables and other structures to hold your data. Now it is time to fill the tables with data. This chapter covers how to insert, update, and delete table data. The chapter after this will finally explain how to extract your long-lost data from the database.

6.1. Inserting Data

When a table is created, it contains no data. The first thing to do before a database can be of much use is to insert data. Data is inserted one row at a time. You can also insert more than one row in a single command, but it is not possible to insert something that is not a complete row. Even if you know only some column values, a complete row must be created.

To create a new row, use the **INSERT** command. The command requires the table name and column values. For example, consider the products table from [Chapter 5](#):

```
CREATE TABLE products (  
    product_no integer,  
    name text,  
    price numeric  
);
```

An example command to insert a row would be:

```
INSERT INTO products VALUES (1, 'Cheese', 9.99);
```

The data values are listed in the order in which the columns appear in the table, separated by commas. Usually, the data values will be literals (constants), but scalar expressions are also allowed.

The above syntax has the drawback that you need to know the order of the columns in the table. To avoid this you can also list the columns explicitly. For example, both of the following commands have the same effect as the one above:

```
INSERT INTO products (product_no, name, price) VALUES (1, 'Cheese', 9.99);  
INSERT INTO products (name, price, product_no) VALUES ('Cheese', 9.99, 1);
```

Many users consider it good practice to always list the column names.

If you don't have values for all the columns, you can omit some of them. In that case, the columns will be filled with their default values. For example:

```
INSERT INTO products (product_no, name) VALUES (1, 'Cheese');  
INSERT INTO products VALUES (1, 'Cheese');
```

The second form is a Postgres Pro extension. It fills the columns from the left with as many values as are given, and the rest will be defaulted.

For clarity, you can also request default values explicitly, for individual columns or for the entire row:

```
INSERT INTO products (product_no, name, price) VALUES (1, 'Cheese', DEFAULT);  
INSERT INTO products DEFAULT VALUES;
```

You can insert multiple rows in a single command:

```
INSERT INTO products (product_no, name, price) VALUES  
    (1, 'Cheese', 9.99),  
    (2, 'Bread', 1.99),  
    (3, 'Milk', 2.99);
```

It is also possible to insert the result of a query (which might be no rows, one row, or many rows):

```
INSERT INTO products (product_no, name, price)  
    SELECT product_no, name, price FROM new_products  
    WHERE release_date = 'today';
```


This provides the full power of the SQL query mechanism ([Chapter 7](#)) for computing the rows to be inserted.

Tip

When inserting a lot of data at the same time, consider using the [COPY](#) command. It is not as flexible as the [INSERT](#) command, but is more efficient. Refer to [Section 14.4](#) for more information on improving bulk loading performance.

6.2. Updating Data

The modification of data that is already in the database is referred to as updating. You can update individual rows, all the rows in a table, or a subset of all rows. Each column can be updated separately; the other columns are not affected.

To update existing rows, use the [UPDATE](#) command. This requires three pieces of information:

1. The name of the table and column to update
2. The new value of the column
3. Which row(s) to update

Recall from [Chapter 5](#) that SQL does not, in general, provide a unique identifier for rows. Therefore it is not always possible to directly specify which row to update. Instead, you specify which conditions a row must meet in order to be updated. Only if you have a primary key in the table (independent of whether you declared it or not) can you reliably address individual rows by choosing a condition that matches the primary key. Graphical database access tools rely on this fact to allow you to update rows individually.

For example, this command updates all products that have a price of 5 to have a price of 10:

```
UPDATE products SET price = 10 WHERE price = 5;
```

This might cause zero, one, or many rows to be updated. It is not an error to attempt an update that does not match any rows.

Let's look at that command in detail. First is the key word `UPDATE` followed by the table name. As usual, the table name can be schema-qualified, otherwise it is looked up in the path. Next is the key word `SET` followed by the column name, an equal sign, and the new column value. The new column value can be any scalar expression, not just a constant. For example, if you want to raise the price of all products by 10% you could use:

```
UPDATE products SET price = price * 1.10;
```

As you see, the expression for the new value can refer to the existing value(s) in the row. We also left out the `WHERE` clause. If it is omitted, it means that all rows in the table are updated. If it is present, only those rows that match the `WHERE` condition are updated. Note that the equals sign in the `SET` clause is an assignment while the one in the `WHERE` clause is a comparison, but this does not create any ambiguity. Of course, the `WHERE` condition does not have to be an equality test. Many other operators are available (see [Chapter 9](#)). But the expression needs to evaluate to a Boolean result.

You can update more than one column in an `UPDATE` command by listing more than one assignment in the `SET` clause. For example:

```
UPDATE mytable SET a = 5, b = 3, c = 1 WHERE a > 0;
```

6.3. Deleting Data

So far we have explained how to add data to tables and how to change data. What remains is to discuss how to remove data that is no longer needed. Just as adding data is only possible in whole rows, you can only remove entire rows from a table. In the previous section we explained that SQL does not provide a way to directly address individual rows. Therefore, removing rows can only be done by specifying

conditions that the rows to be removed have to match. If you have a primary key in the table then you can specify the exact row. But you can also remove groups of rows matching a condition, or you can remove all rows in the table at once.

You use the **DELETE** command to remove rows; the syntax is very similar to the **UPDATE** command. For instance, to remove all rows from the products table that have a price of 10, use:

```
DELETE FROM products WHERE price = 10;
```

If you simply write:

```
DELETE FROM products;
```

then all rows in the table will be deleted! Caveat programmer.

6.4. Returning Data from Modified Rows

Sometimes it is useful to obtain data from modified rows while they are being manipulated. The **INSERT**, **UPDATE**, **DELETE**, and **MERGE** commands all have an optional **RETURNING** clause that supports this. Use of **RETURNING** avoids performing an extra database query to collect the data, and is especially valuable when it would otherwise be difficult to identify the modified rows reliably.

The allowed contents of a **RETURNING** clause are the same as a **SELECT** command's output list (see [Section 7.3](#)). It can contain column names of the command's target table, or value expressions using those columns. A common shorthand is **RETURNING ***, which selects all columns of the target table in order.

In an **INSERT**, the data available to **RETURNING** is the row as it was inserted. This is not so useful in trivial inserts, since it would just repeat the data provided by the client. But it can be very handy when relying on computed default values. For example, when using a **serial** column to provide unique identifiers, **RETURNING** can return the ID assigned to a new row:

```
CREATE TABLE users (firstname text, lastname text, id serial primary key);
```

```
INSERT INTO users (firstname, lastname) VALUES ('Joe', 'Cool') RETURNING id;
```

The **RETURNING** clause is also very useful with **INSERT ... SELECT**.

In an **UPDATE**, the data available to **RETURNING** is the new content of the modified row. For example:

```
UPDATE products SET price = price * 1.10
WHERE price <= 99.99
RETURNING name, price AS new_price;
```

In a **DELETE**, the data available to **RETURNING** is the content of the deleted row. For example:

```
DELETE FROM products
WHERE obsolescence_date = 'today'
RETURNING *;
```

In a **MERGE**, the data available to **RETURNING** is the content of the source row plus the content of the inserted, updated, or deleted target row. Since it is quite common for the source and target to have many of the same columns, specifying **RETURNING *** can lead to a lot of duplicated columns, so it is often more useful to qualify it so as to return just the source or target row. For example:

```
MERGE INTO products p USING new_products n ON p.product_no = n.product_no
WHEN NOT MATCHED THEN INSERT VALUES (n.product_no, n.name, n.price)
WHEN MATCHED THEN UPDATE SET name = n.name, price = n.price
RETURNING p.*;
```

If there are triggers ([Chapter 40](#)) on the target table, the data available to **RETURNING** is the row as modified by the triggers. Thus, inspecting columns computed by triggers is another common use-case for **RETURNING**.

Chapter 7. Queries

The previous chapters explained how to create tables, how to fill them with data, and how to manipulate that data. Now we finally discuss how to retrieve the data from the database.

7.1. Overview

The process of retrieving or the command to retrieve data from a database is called a *query*. In SQL the `SELECT` command is used to specify queries. The general syntax of the `SELECT` command is

```
[WITH with_queries] SELECT select_list FROM table_expression [sort_specification]
```

The following sections describe the details of the select list, the table expression, and the sort specification. `WITH` queries are treated last since they are an advanced feature.

A simple kind of query has the form:

```
SELECT * FROM table1;
```

Assuming that there is a table called `table1`, this command would retrieve all rows and all user-defined columns from `table1`. (The method of retrieval depends on the client application. For example, the `psql` program will display an ASCII-art table on the screen, while client libraries will offer functions to extract individual values from the query result.) The select list specification `*` means all columns that the table expression happens to provide. A select list can also select a subset of the available columns or make calculations using the columns. For example, if `table1` has columns named `a`, `b`, and `c` (and perhaps others) you can make the following query:

```
SELECT a, b + c FROM table1;
```

(assuming that `b` and `c` are of a numerical data type). See [Section 7.3](#) for more details.

`FROM table1` is a simple kind of table expression: it reads just one table. In general, table expressions can be complex constructs of base tables, joins, and subqueries. But you can also omit the table expression entirely and use the `SELECT` command as a calculator:

```
SELECT 3 * 4;
```

This is more useful if the expressions in the select list return varying results. For example, you could call a function this way:

```
SELECT random();
```

7.2. Table Expressions

A *table expression* computes a table. The table expression contains a `FROM` clause that is optionally followed by `WHERE`, `GROUP BY`, and `HAVING` clauses. Trivial table expressions simply refer to a table on disk, a so-called base table, but more complex expressions can be used to modify or combine base tables in various ways.

The optional `WHERE`, `GROUP BY`, and `HAVING` clauses in the table expression specify a pipeline of successive transformations performed on the table derived in the `FROM` clause. All these transformations produce a virtual table that provides the rows that are passed to the select list to compute the output rows of the query.

7.2.1. The `FROM` Clause

The `FROM` clause derives a table from one or more other tables given in a comma-separated table reference list.

```
FROM table_reference [, table_reference [, ...]]
```

A table reference can be a table name (possibly schema-qualified), or a derived table such as a subquery, a `JOIN` construct, or complex combinations of these. If more than one table reference is listed in the `FROM` clause, the tables are cross-joined (that is, the Cartesian product of their rows is formed; see below).

The result of the `FROM` list is an intermediate virtual table that can then be subject to transformations by the `WHERE`, `GROUP BY`, and `HAVING` clauses and is finally the result of the overall table expression.

When a table reference names a table that is the parent of a table inheritance hierarchy, the table reference produces rows of not only that table but all of its descendant tables, unless the key word `ONLY` precedes the table name. However, the reference produces only the columns that appear in the named table — any columns added in subtables are ignored.

Instead of writing `ONLY` before the table name, you can write `*` after the table name to explicitly specify that descendant tables are included. There is no real reason to use this syntax any more, because searching descendant tables is now always the default behavior. However, it is supported for compatibility with older releases.

7.2.1.1. Joined Tables

A joined table is a table derived from two other (real or derived) tables according to the rules of the particular join type. Inner, outer, and cross-joins are available. The general syntax of a joined table is

```
T1 join_type T2 [ join_condition ]
```

Joins of all types can be chained together, or nested: either or both *T1* and *T2* can be joined tables. Parentheses can be used around `JOIN` clauses to control the join order. In the absence of parentheses, `JOIN` clauses nest left-to-right.

Join Types

Cross join

```
T1 CROSS JOIN T2
```

For every possible combination of rows from *T1* and *T2* (i.e., a Cartesian product), the joined table will contain a row consisting of all columns in *T1* followed by all columns in *T2*. If the tables have *N* and *M* rows respectively, the joined table will have *N * M* rows.

`FROM T1 CROSS JOIN T2` is equivalent to `FROM T1 INNER JOIN T2 ON TRUE` (see below). It is also equivalent to `FROM T1, T2`.

Note

This latter equivalence does not hold exactly when more than two tables appear, because `JOIN` binds more tightly than comma. For example `FROM T1 CROSS JOIN T2 INNER JOIN T3 ON condition` is not the same as `FROM T1, T2 INNER JOIN T3 ON condition` because the *condition* can reference *T1* in the first case but not the second.

Qualified joins

```
T1 { [INNER] | { LEFT | RIGHT | FULL } [OUTER] } JOIN T2 ON boolean_expression
T1 { [INNER] | { LEFT | RIGHT | FULL } [OUTER] } JOIN T2 USING ( join column list )
T1 NATURAL { [INNER] | { LEFT | RIGHT | FULL } [OUTER] } JOIN T2
```

The words `INNER` and `OUTER` are optional in all forms. `INNER` is the default; `LEFT`, `RIGHT`, and `FULL` imply an outer join.

The *join condition* is specified in the `ON` or `USING` clause, or implicitly by the word `NATURAL`. The join condition determines which rows from the two source tables are considered to “match”, as explained in detail below.

The possible types of qualified join are:

```
INNER JOIN
```

For each row *R1* of *T1*, the joined table has a row for each row in *T2* that satisfies the join condition with *R1*.

LEFT OUTER JOIN

First, an inner join is performed. Then, for each row in *T1* that does not satisfy the join condition with any row in *T2*, a joined row is added with null values in columns of *T2*. Thus, the joined table always has at least one row for each row in *T1*.

RIGHT OUTER JOIN

First, an inner join is performed. Then, for each row in *T2* that does not satisfy the join condition with any row in *T1*, a joined row is added with null values in columns of *T1*. This is the converse of a left join: the result table will always have a row for each row in *T2*.

FULL OUTER JOIN

First, an inner join is performed. Then, for each row in *T1* that does not satisfy the join condition with any row in *T2*, a joined row is added with null values in columns of *T2*. Also, for each row of *T2* that does not satisfy the join condition with any row in *T1*, a joined row with null values in the columns of *T1* is added.

The `ON` clause is the most general kind of join condition: it takes a Boolean value expression of the same kind as is used in a `WHERE` clause. A pair of rows from *T1* and *T2* match if the `ON` expression evaluates to true.

The `USING` clause is a shorthand that allows you to take advantage of the specific situation where both sides of the join use the same name for the joining column(s). It takes a comma-separated list of the shared column names and forms a join condition that includes an equality comparison for each one. For example, joining *T1* and *T2* with `USING (a, b)` produces the join condition `ON T1.a = T2.a AND T1.b = T2.b`.

Furthermore, the output of `JOIN USING` suppresses redundant columns: there is no need to print both of the matched columns, since they must have equal values. While `JOIN ON` produces all columns from *T1* followed by all columns from *T2*, `JOIN USING` produces one output column for each of the listed column pairs (in the listed order), followed by any remaining columns from *T1*, followed by any remaining columns from *T2*.

Finally, `NATURAL` is a shorthand form of `USING`: it forms a `USING` list consisting of all column names that appear in both input tables. As with `USING`, these columns appear only once in the output table. If there are no common column names, `NATURAL JOIN` behaves like `CROSS JOIN`.

Note

`USING` is reasonably safe from column changes in the joined relations since only the listed columns are combined. `NATURAL` is considerably more risky since any schema changes to either relation that cause a new matching column name to be present will cause the join to combine that new column as well.

To put this together, assume we have tables *t1*:

num	name
1	a
2	b
3	c

and *t2*:

num	value
1	xxx
3	yyy

5 | zzz

then we get the following results for the various joins:

=> **SELECT * FROM t1 CROSS JOIN t2;**

num	name	num	value
1	a	1	xxx
1	a	3	yyy
1	a	5	zzz
2	b	1	xxx
2	b	3	yyy
2	b	5	zzz
3	c	1	xxx
3	c	3	yyy
3	c	5	zzz

(9 rows)

=> **SELECT * FROM t1 INNER JOIN t2 ON t1.num = t2.num;**

num	name	num	value
1	a	1	xxx
3	c	3	yyy

(2 rows)

=> **SELECT * FROM t1 INNER JOIN t2 USING (num);**

num	name	value
1	a	xxx
3	c	yyy

(2 rows)

=> **SELECT * FROM t1 NATURAL INNER JOIN t2;**

num	name	value
1	a	xxx
3	c	yyy

(2 rows)

=> **SELECT * FROM t1 LEFT JOIN t2 ON t1.num = t2.num;**

num	name	num	value
1	a	1	xxx
2	b		
3	c	3	yyy

(3 rows)

=> **SELECT * FROM t1 LEFT JOIN t2 USING (num);**

num	name	value
1	a	xxx
2	b	
3	c	yyy

(3 rows)

=> **SELECT * FROM t1 RIGHT JOIN t2 ON t1.num = t2.num;**

num	name	num	value
1	a	1	xxx

```

3 | c      | 3 | yyy
  |         | 5 | zzz
(3 rows)

```

```

=> SELECT * FROM t1 FULL JOIN t2 ON t1.num = t2.num;
 num | name | num | value
-----+-----+-----+-----
  1  | a    |  1  | xxx
  2  | b    |     |
  3  | c    |  3  | yyy
     |     |  5  | zzz
(4 rows)

```

The join condition specified with `ON` can also contain conditions that do not relate directly to the join. This can prove useful for some queries but needs to be thought out carefully. For example:

```

=> SELECT * FROM t1 LEFT JOIN t2 ON t1.num = t2.num AND t2.value = 'xxx';
 num | name | num | value
-----+-----+-----+-----
  1  | a    |  1  | xxx
  2  | b    |     |
  3  | c    |     |
(3 rows)

```

Notice that placing the restriction in the `WHERE` clause produces a different result:

```

=> SELECT * FROM t1 LEFT JOIN t2 ON t1.num = t2.num WHERE t2.value = 'xxx';
 num | name | num | value
-----+-----+-----+-----
  1  | a    |  1  | xxx
(1 row)

```

This is because a restriction placed in the `ON` clause is processed *before* the join, while a restriction placed in the `WHERE` clause is processed *after* the join. That does not matter with inner joins, but it matters a lot with outer joins.

7.2.1.2. Table and Column Aliases

A temporary name can be given to tables and complex table references to be used for references to the derived table in the rest of the query. This is called a *table alias*.

To create a table alias, write

```
FROM table_reference AS alias
```

or

```
FROM table_reference alias
```

The `AS` key word is optional noise. *alias* can be any identifier.

A typical application of table aliases is to assign short identifiers to long table names to keep the join clauses readable. For example:

```
SELECT * FROM some_very_long_table_name s JOIN another_fairly_long_name a ON s.id =
a.num;
```

The alias becomes the new name of the table reference so far as the current query is concerned — it is not allowed to refer to the table by the original name elsewhere in the query. Thus, this is not valid:

```
SELECT * FROM my_table AS m WHERE my_table.a > 5;      -- wrong
```

Table aliases are mainly for notational convenience, but it is necessary to use them when joining a table to itself, e.g.:

```
SELECT * FROM people AS mother JOIN people AS child ON mother.id = child.mother_id;
```

Parentheses are used to resolve ambiguities. In the following example, the first statement assigns the alias `b` to the second instance of `my_table`, but the second statement assigns the alias to the result of the join:

```
SELECT * FROM my_table AS a CROSS JOIN my_table AS b ...
SELECT * FROM (my_table AS a CROSS JOIN my_table) AS b ...
```

Another form of table aliasing gives temporary names to the columns of the table, as well as the table itself:

```
FROM table_reference [AS] alias ( column1 [, column2 [, ...]] )
```

If fewer column aliases are specified than the actual table has columns, the remaining columns are not renamed. This syntax is especially useful for self-joins or subqueries.

When an alias is applied to the output of a `JOIN` clause, the alias hides the original name(s) within the `JOIN`. For example:

```
SELECT a.* FROM my_table AS a JOIN your_table AS b ON ...
```

is valid SQL, but:

```
SELECT a.* FROM (my_table AS a JOIN your_table AS b ON ...) AS c
```

is not valid; the table alias `a` is not visible outside the alias `c`.

7.2.1.3. Subqueries

Subqueries specifying a derived table must be enclosed in parentheses. They may be assigned a table alias name, and optionally column alias names (as in [Section 7.2.1.2](#)). For example:

```
FROM (SELECT * FROM table1) AS alias_name
```

This example is equivalent to `FROM table1 AS alias_name`. More interesting cases, which cannot be reduced to a plain join, arise when the subquery involves grouping or aggregation.

A subquery can also be a `VALUES` list:

```
FROM (VALUES ('anne', 'smith'), ('bob', 'jones'), ('joe', 'blow'))
     AS names(first, last)
```

Again, a table alias is optional. Assigning alias names to the columns of the `VALUES` list is optional, but is good practice. For more information see [Section 7.7](#).

According to the SQL standard, a table alias name must be supplied for a subquery. Postgres Pro allows `AS` and the alias to be omitted, but writing one is good practice in SQL code that might be ported to another system.

7.2.1.4. Table Functions

Table functions are functions that produce a set of rows, made up of either base data types (scalar types) or composite data types (table rows). They are used like a table, view, or subquery in the `FROM` clause of a query. Columns returned by table functions can be included in `SELECT`, `JOIN`, or `WHERE` clauses in the same manner as columns of a table, view, or subquery.

Table functions may also be combined using the `ROWS FROM` syntax, with the results returned in parallel columns; the number of result rows in this case is that of the largest function result, with smaller results padded with null values to match.

```
function_call [WITH ORDINALITY] [[AS] table_alias [(column_alias [, ... ])] ]
ROWS FROM( function_call [, ... ] ) [WITH ORDINALITY] [[AS] table_alias [(column_alias
[, ... ])] ]
```


If the `WITH ORDINALITY` clause is specified, an additional column of type `bigint` will be added to the function result columns. This column numbers the rows of the function result set, starting from 1. (This is a generalization of the SQL-standard syntax for `UNNEST ... WITH ORDINALITY`.) By default, the ordinal column is called `ordinality`, but a different column name can be assigned to it using an `AS` clause.

The special table function `UNNEST` may be called with any number of array parameters, and it returns a corresponding number of columns, as if `UNNEST` ([Section 9.19](#)) had been called on each parameter separately and combined using the `ROWS FROM` construct.

```
UNNEST( array_expression [, ... ] ) [WITH ORDINALITY] [[AS] table_alias [(column_alias
[, ... ])]]
```

If no `table_alias` is specified, the function name is used as the table name; in the case of a `ROWS FROM()` construct, the first function's name is used.

If column aliases are not supplied, then for a function returning a base data type, the column name is also the same as the function name. For a function returning a composite type, the result columns get the names of the individual attributes of the type.

Some examples:

```
CREATE TABLE foo (fooid int, foosubid int, fooname text);
```

```
CREATE FUNCTION getfoo(int) RETURNS SETOF foo AS $$
    SELECT * FROM foo WHERE fooid = $1;
$$ LANGUAGE SQL;
```

```
SELECT * FROM getfoo(1) AS t1;
```

```
SELECT * FROM foo
    WHERE foosubid IN (
                        SELECT foosubid
                        FROM getfoo(foo.fooid) z
                        WHERE z.fooid = foo.fooid
                    );
```

```
CREATE VIEW vw_getfoo AS SELECT * FROM getfoo(1);
```

```
SELECT * FROM vw_getfoo;
```

In some cases it is useful to define table functions that can return different column sets depending on how they are invoked. To support this, the table function can be declared as returning the pseudo-type `record` with no `OUT` parameters. When such a function is used in a query, the expected row structure must be specified in the query itself, so that the system can know how to parse and plan the query. This syntax looks like:

```
function_call [AS] alias (column_definition [, ... ])
function_call AS [alias] (column_definition [, ... ])
ROWS FROM( ... function_call AS (column_definition [, ... ]) [, ... ] )
```

When not using the `ROWS FROM()` syntax, the `column_definition` list replaces the column alias list that could otherwise be attached to the `FROM` item; the names in the column definitions serve as column aliases. When using the `ROWS FROM()` syntax, a `column_definition` list can be attached to each member function separately; or if there is only one member function and no `WITH ORDINALITY` clause, a `column_definition` list can be written in place of a column alias list following `ROWS FROM()`.

Consider this example:

```
SELECT *
    FROM dblink('dbname=mydb', 'SELECT proname, prosrc FROM pg_proc')
    AS t1(proname name, prosrc text)
```

```
WHERE pronom LIKE 'bytea%';
```

The `dblink` function (part of the `dblink` module) executes a remote query. It is declared to return `record` since it might be used for any kind of query. The actual column set must be specified in the calling query so that the parser knows, for example, what `*` should expand to.

This example uses `ROWS FROM`:

```
SELECT *
FROM ROWS FROM
(
    json_to_recordset('[{ "a":40, "b":"foo"}, {"a":100, "b":"bar"}]')
    AS (a INTEGER, b TEXT),
    generate_series(1, 3)
) AS x (p, q, s)
ORDER BY p;
```

p	q	s
40	foo	1
100	bar	2
		3

It joins two functions into a single `FROM` target. `json_to_recordset()` is instructed to return two columns, the first integer and the second text. The result of `generate_series()` is used directly. The `ORDER BY` clause sorts the column values as integers.

7.2.1.5. LATERAL Subqueries

Subqueries appearing in `FROM` can be preceded by the key word `LATERAL`. This allows them to reference columns provided by preceding `FROM` items. (Without `LATERAL`, each subquery is evaluated independently and so cannot cross-reference any other `FROM` item.)

Table functions appearing in `FROM` can also be preceded by the key word `LATERAL`, but for functions the key word is optional; the function's arguments can contain references to columns provided by preceding `FROM` items in any case.

A `LATERAL` item can appear at the top level in the `FROM` list, or within a `JOIN` tree. In the latter case it can also refer to any items that are on the left-hand side of a `JOIN` that it is on the right-hand side of.

When a `FROM` item contains `LATERAL` cross-references, evaluation proceeds as follows: for each row of the `FROM` item providing the cross-referenced column(s), or set of rows of multiple `FROM` items providing the columns, the `LATERAL` item is evaluated using that row or row set's values of the columns. The resulting row(s) are joined as usual with the rows they were computed from. This is repeated for each row or set of rows from the column source table(s).

A trivial example of `LATERAL` is

```
SELECT * FROM foo, LATERAL (SELECT * FROM bar WHERE bar.id = foo.bar_id) ss;
```

This is not especially useful since it has exactly the same result as the more conventional

```
SELECT * FROM foo, bar WHERE bar.id = foo.bar_id;
```

`LATERAL` is primarily useful when the cross-referenced column is necessary for computing the row(s) to be joined. A common application is providing an argument value for a set-returning function. For example, supposing that `vertices(polygon)` returns the set of vertices of a polygon, we could identify close-together vertices of polygons stored in a table with:

```
SELECT p1.id, p2.id, v1, v2
FROM polygons p1, polygons p2,
    LATERAL vertices(p1.poly) v1,
    LATERAL vertices(p2.poly) v2
```

```
WHERE (v1 <-> v2) < 10 AND p1.id != p2.id;
```

This query could also be written

```
SELECT p1.id, p2.id, v1, v2
FROM polygons p1 CROSS JOIN LATERAL vertices(p1.poly) v1,
     polygons p2 CROSS JOIN LATERAL vertices(p2.poly) v2
WHERE (v1 <-> v2) < 10 AND p1.id != p2.id;
```

or in several other equivalent formulations. (As already mentioned, the `LATERAL` key word is unnecessary in this example, but we use it for clarity.)

It is often particularly handy to `LEFT JOIN` to a `LATERAL` subquery, so that source rows will appear in the result even if the `LATERAL` subquery produces no rows for them. For example, if `get_product_names()` returns the names of products made by a manufacturer, but some manufacturers in our table currently produce no products, we could find out which ones those are like this:

```
SELECT m.name
FROM manufacturers m LEFT JOIN LATERAL get_product_names(m.id) pname ON true
WHERE pname IS NULL;
```

7.2.2. The WHERE Clause

The syntax of the `WHERE` clause is

```
WHERE search_condition
```

where *search_condition* is any value expression (see [Section 4.2](#)) that returns a value of type `boolean`.

After the processing of the `FROM` clause is done, each row of the derived virtual table is checked against the search condition. If the result of the condition is true, the row is kept in the output table, otherwise (i.e., if the result is false or null) it is discarded. The search condition typically references at least one column of the table generated in the `FROM` clause; this is not required, but otherwise the `WHERE` clause will be fairly useless.

Note

The join condition of an inner join can be written either in the `WHERE` clause or in the `JOIN` clause. For example, these table expressions are equivalent:

```
FROM a, b WHERE a.id = b.id AND b.val > 5
```

and:

```
FROM a INNER JOIN b ON (a.id = b.id) WHERE b.val > 5
```

or perhaps even:

```
FROM a NATURAL JOIN b WHERE b.val > 5
```

Which one of these you use is mainly a matter of style. The `JOIN` syntax in the `FROM` clause is probably not as portable to other SQL database management systems, even though it is in the SQL standard. For outer joins there is no choice: they must be done in the `FROM` clause. The `ON` or `USING` clause of an outer join is *not* equivalent to a `WHERE` condition, because it results in the addition of rows (for unmatched input rows) as well as the removal of rows in the final result.

Here are some examples of `WHERE` clauses:

```
SELECT ... FROM fdt WHERE c1 > 5
```

```
SELECT ... FROM fdt WHERE c1 IN (1, 2, 3)
```

```
SELECT ... FROM fdt WHERE c1 IN (SELECT c1 FROM t2)
```

```
SELECT ... FROM fdt WHERE c1 IN (SELECT c3 FROM t2 WHERE c2 = fdt.c1 + 10)
```

```
SELECT ... FROM fdt WHERE c1 BETWEEN (SELECT c3 FROM t2 WHERE c2 = fdt.c1 + 10) AND 100
```

```
SELECT ... FROM fdt WHERE EXISTS (SELECT c1 FROM t2 WHERE c2 > fdt.c1)
```

`fdt` is the table derived in the `FROM` clause. Rows that do not meet the search condition of the `WHERE` clause are eliminated from `fdt`. Notice the use of scalar subqueries as value expressions. Just like any other query, the subqueries can employ complex table expressions. Notice also how `fdt` is referenced in the subqueries. Qualifying `c1` as `fdt.c1` is only necessary if `c1` is also the name of a column in the derived input table of the subquery. But qualifying the column name adds clarity even when it is not needed. This example shows how the column naming scope of an outer query extends into its inner queries.

7.2.3. The `GROUP BY` and `HAVING` Clauses

After passing the `WHERE` filter, the derived input table might be subject to grouping, using the `GROUP BY` clause, and elimination of group rows using the `HAVING` clause.

```
SELECT select_list
  FROM ...
  [WHERE ...]
  GROUP BY grouping_column_reference [, grouping_column_reference]...
```

The `GROUP BY` clause is used to group together those rows in a table that have the same values in all the columns listed. The order in which the columns are listed does not matter. The effect is to combine each set of rows having common values into one group row that represents all rows in the group. This is done to eliminate redundancy in the output and/or compute aggregates that apply to these groups. For instance:

```
=> SELECT * FROM test1;
 x | y
---+---
 a | 3
 c | 2
 b | 5
 a | 1
(4 rows)
```

```
=> SELECT x FROM test1 GROUP BY x;
 x
---
 a
 b
 c
(3 rows)
```

In the second query, we could not have written `SELECT * FROM test1 GROUP BY x`, because there is no single value for the column `y` that could be associated with each group. The grouped-by columns can be referenced in the select list since they have a single value in each group.

In general, if a table is grouped, columns that are not listed in `GROUP BY` cannot be referenced except in aggregate expressions. An example with aggregate expressions is:

```
=> SELECT x, sum(y) FROM test1 GROUP BY x;
 x | sum
---+-----
 a |    4
 b |    5
 c |    2
(3 rows)
```

Here `sum` is an aggregate function that computes a single value over the entire group. More information about the available aggregate functions can be found in [Section 9.21](#).

Tip

Grouping without aggregate expressions effectively calculates the set of distinct values in a column. This can also be achieved using the `DISTINCT` clause (see [Section 7.3.3](#)).

Here is another example: it calculates the total sales for each product (rather than the total sales of all products):

```
SELECT product_id, p.name, (sum(s.units) * p.price) AS sales
FROM products p LEFT JOIN sales s USING (product_id)
GROUP BY product_id, p.name, p.price;
```

In this example, the columns `product_id`, `p.name`, and `p.price` must be in the `GROUP BY` clause since they are referenced in the query select list (but see below). The column `s.units` does not have to be in the `GROUP BY` list since it is only used in an aggregate expression (`sum(...)`), which represents the sales of a product. For each product, the query returns a summary row about all sales of the product.

If the products table is set up so that, say, `product_id` is the primary key, then it would be enough to group by `product_id` in the above example, since name and price would be *functionally dependent* on the product ID, and so there would be no ambiguity about which name and price value to return for each product ID group.

In strict SQL, `GROUP BY` can only group by columns of the source table but Postgres Pro extends this to also allow `GROUP BY` to group by columns in the select list. Grouping by value expressions instead of simple column names is also allowed.

If a table has been grouped using `GROUP BY`, but only certain groups are of interest, the `HAVING` clause can be used, much like a `WHERE` clause, to eliminate groups from the result. The syntax is:

```
SELECT select_list FROM ... [WHERE ...] GROUP BY ... HAVING boolean_expression
```

Expressions in the `HAVING` clause can refer both to grouped expressions and to ungrouped expressions (which necessarily involve an aggregate function).

Example:

```
=> SELECT x, sum(y) FROM test1 GROUP BY x HAVING sum(y) > 3;
```

```
x | sum
---+-----
a |    4
b |    5
(2 rows)
```

```
=> SELECT x, sum(y) FROM test1 GROUP BY x HAVING x < 'c';
```

```
x | sum
---+-----
a |    4
b |    5
(2 rows)
```

Again, a more realistic example:

```
SELECT product_id, p.name, (sum(s.units) * (p.price - p.cost)) AS profit
FROM products p LEFT JOIN sales s USING (product_id)
WHERE s.date > CURRENT_DATE - INTERVAL '4 weeks'
GROUP BY product_id, p.name, p.price, p.cost
```

```
HAVING sum(p.price * s.units) > 5000;
```

In the example above, the `WHERE` clause is selecting rows by a column that is not grouped (the expression is only true for sales during the last four weeks), while the `HAVING` clause restricts the output to groups with total gross sales over 5000. Note that the aggregate expressions do not necessarily need to be the same in all parts of the query.

If a query contains aggregate function calls, but no `GROUP BY` clause, grouping still occurs: the result is a single group row (or perhaps no rows at all, if the single row is then eliminated by `HAVING`). The same is true if it contains a `HAVING` clause, even without any aggregate function calls or `GROUP BY` clause.

7.2.4. GROUPING SETS, CUBE, and ROLLUP

More complex grouping operations than those described above are possible using the concept of *grouping sets*. The data selected by the `FROM` and `WHERE` clauses is grouped separately by each specified grouping set, aggregates computed for each group just as for simple `GROUP BY` clauses, and then the results returned. For example:

```
=> SELECT * FROM items_sold;
```

brand	size	sales
Foo	L	10
Foo	M	20
Bar	M	15
Bar	L	5

(4 rows)

```
=> SELECT brand, size, sum(sales) FROM items_sold GROUP BY GROUPING SETS ((brand),
    (size), ());
```

brand	size	sum
Foo		30
Bar		20
	L	15
	M	35
		50

(5 rows)

Each sublist of `GROUPING SETS` may specify zero or more columns or expressions and is interpreted the same way as though it were directly in the `GROUP BY` clause. An empty grouping set means that all rows are aggregated down to a single group (which is output even if no input rows were present), as described above for the case of aggregate functions with no `GROUP BY` clause.

References to the grouping columns or expressions are replaced by null values in result rows for grouping sets in which those columns do not appear. To distinguish which grouping a particular output row resulted from, see [Table 9.64](#).

A shorthand notation is provided for specifying two common types of grouping set. A clause of the form

```
ROLLUP ( e1, e2, e3, ... )
```

represents the given list of expressions and all prefixes of the list including the empty list; thus it is equivalent to

```
GROUPING SETS (
    ( e1, e2, e3, ... ),
    ...
    ( e1, e2 ),
    ( e1 ),
    ( )
)
```

This is commonly used for analysis over hierarchical data; e.g., total salary by department, division, and company-wide total.

A clause of the form

```
CUBE ( e1, e2, ... )
```

represents the given list and all of its possible subsets (i.e., the power set). Thus

```
CUBE ( a, b, c )
```

is equivalent to

```
GROUPING SETS (
  ( a, b, c ),
  ( a, b   ),
  ( a,    c ),
  ( a     ),
  (    b, c ),
  (    b   ),
  (     c ),
  (      )
)
```

The individual elements of a `CUBE` or `ROLLUP` clause may be either individual expressions, or sublists of elements in parentheses. In the latter case, the sublists are treated as single units for the purposes of generating the individual grouping sets. For example:

```
CUBE ( (a, b), (c, d) )
```

is equivalent to

```
GROUPING SETS (
  ( a, b, c, d ),
  ( a, b       ),
  (      c, d ),
  (           )
)
```

and

```
ROLLUP ( a, (b, c), d )
```

is equivalent to

```
GROUPING SETS (
  ( a, b, c, d ),
  ( a, b, c     ),
  ( a           ),
  (            )
)
```

The `CUBE` and `ROLLUP` constructs can be used either directly in the `GROUP BY` clause, or nested inside a `GROUPING SETS` clause. If one `GROUPING SETS` clause is nested inside another, the effect is the same as if all the elements of the inner clause had been written directly in the outer clause.

If multiple grouping items are specified in a single `GROUP BY` clause, then the final list of grouping sets is the Cartesian product of the individual items. For example:

```
GROUP BY a, CUBE (b, c), GROUPING SETS ((d), (e))
```

is equivalent to

```
GROUP BY GROUPING SETS (
  (a, b, c, d), (a, b, c, e),
  (a, b, d),   (a, b, e),

```

```
(a, c, d),    (a, c, e),
(a, d),      (a, e)
)
```

When specifying multiple grouping items together, the final set of grouping sets might contain duplicates. For example:

```
GROUP BY ROLLUP (a, b), ROLLUP (a, c)
```

is equivalent to

```
GROUP BY GROUPING SETS (
    (a, b, c),
    (a, b),
    (a, b),
    (a, c),
    (a),
    (a),
    (a, c),
    (a),
    ()
)
```

If these duplicates are undesirable, they can be removed using the `DISTINCT` clause directly on the `GROUP BY`. Therefore:

```
GROUP BY DISTINCT ROLLUP (a, b), ROLLUP (a, c)
```

is equivalent to

```
GROUP BY GROUPING SETS (
    (a, b, c),
    (a, b),
    (a, c),
    (a),
    ()
)
```

This is not the same as using `SELECT DISTINCT` because the output rows may still contain duplicates. If any of the ungrouped columns contains `NULL`, it will be indistinguishable from the `NULL` used when that same column is grouped.

Note

The construct `(a, b)` is normally recognized in expressions as a [row constructor](#). Within the `GROUP BY` clause, this does not apply at the top levels of expressions, and `(a, b)` is parsed as a list of expressions as described above. If for some reason you *need* a row constructor in a grouping expression, use `ROW(a, b)`.

7.2.5. Window Function Processing

If the query contains any window functions (see [Section 3.5](#), [Section 9.22](#) and [Section 4.2.8](#)), these functions are evaluated after any grouping, aggregation, and `HAVING` filtering is performed. That is, if the query uses any aggregates, `GROUP BY`, or `HAVING`, then the rows seen by the window functions are the group rows instead of the original table rows from `FROM/WHERE`.

When multiple window functions are used, all the window functions having equivalent `PARTITION BY` and `ORDER BY` clauses in their window definitions are guaranteed to see the same ordering of the input rows, even if the `ORDER BY` does not uniquely determine the ordering. However, no guarantees are made about the evaluation of functions having different `PARTITION BY` or `ORDER BY` specifications. (In such

cases a sort step is typically required between the passes of window function evaluations, and the sort is not guaranteed to preserve ordering of rows that its `ORDER BY` sees as equivalent.)

Currently, window functions always require presorted data, and so the query output will be ordered according to one or another of the window functions' `PARTITION BY/ORDER BY` clauses. It is not recommended to rely on this, however. Use an explicit top-level `ORDER BY` clause if you want to be sure the results are sorted in a particular way.

7.3. Select Lists

As shown in the previous section, the table expression in the `SELECT` command constructs an intermediate virtual table by possibly combining tables, views, eliminating rows, grouping, etc. This table is finally passed on to processing by the *select list*. The select list determines which *columns* of the intermediate table are actually output.

7.3.1. Select-List Items

The simplest kind of select list is `*` which emits all columns that the table expression produces. Otherwise, a select list is a comma-separated list of value expressions (as defined in [Section 4.2](#)). For instance, it could be a list of column names:

```
SELECT a, b, c FROM ...
```

The columns names `a`, `b`, and `c` are either the actual names of the columns of tables referenced in the `FROM` clause, or the aliases given to them as explained in [Section 7.2.1.2](#). The name space available in the select list is the same as in the `WHERE` clause, unless grouping is used, in which case it is the same as in the `HAVING` clause.

If more than one table has a column of the same name, the table name must also be given, as in:

```
SELECT tbl1.a, tbl2.a, tbl1.b FROM ...
```

When working with multiple tables, it can also be useful to ask for all the columns of a particular table:

```
SELECT tbl1.*, tbl2.a FROM ...
```

See [Section 8.16.5](#) for more about the `table_name.*` notation.

If an arbitrary value expression is used in the select list, it conceptually adds a new virtual column to the returned table. The value expression is evaluated once for each result row, with the row's values substituted for any column references. But the expressions in the select list do not have to reference any columns in the table expression of the `FROM` clause; they can be constant arithmetic expressions, for instance.

7.3.2. Column Labels

The entries in the select list can be assigned names for subsequent processing, such as for use in an `ORDER BY` clause or for display by the client application. For example:

```
SELECT a AS value, b + c AS sum FROM ...
```

If no output column name is specified using `AS`, the system assigns a default column name. For simple column references, this is the name of the referenced column. For function calls, this is the name of the function. For complex expressions, the system will generate a generic name.

The `AS` key word is usually optional, but in some cases where the desired column name matches a Postgres Pro key word, you must write `AS` or double-quote the column name in order to avoid ambiguity. ([Appendix C](#) shows which key words require `AS` to be used as a column label.) For example, `FROM` is one such key word, so this does not work:

```
SELECT a from, b + c AS sum FROM ...
```

but either of these do:

```
SELECT a AS from, b + c AS sum FROM ...
```

```
SELECT a "from", b + c AS sum FROM ...
```

For greatest safety against possible future key word additions, it is recommended that you always either write `AS` or double-quote the output column name.

Note

The naming of output columns here is different from that done in the `FROM` clause (see [Section 7.2.1.2](#)). It is possible to rename the same column twice, but the name assigned in the select list is the one that will be passed on.

7.3.3. DISTINCT

After the select list has been processed, the result table can optionally be subject to the elimination of duplicate rows. The `DISTINCT` key word is written directly after `SELECT` to specify this:

```
SELECT DISTINCT select_list ...
```

(Instead of `DISTINCT` the key word `ALL` can be used to specify the default behavior of retaining all rows.)

Obviously, two rows are considered distinct if they differ in at least one column value. Null values are considered equal in this comparison.

Alternatively, an arbitrary expression can determine what rows are to be considered distinct:

```
SELECT DISTINCT ON (expression [, expression ...]) select_list ...
```

Here *expression* is an arbitrary value expression that is evaluated for all rows. A set of rows for which all the expressions are equal are considered duplicates, and only the first row of the set is kept in the output. Note that the “first row” of a set is unpredictable unless the query is sorted on enough columns to guarantee a unique ordering of the rows arriving at the `DISTINCT` filter. (`DISTINCT ON` processing occurs after `ORDER BY` sorting.)

The `DISTINCT ON` clause is not part of the SQL standard and is sometimes considered bad style because of the potentially indeterminate nature of its results. With judicious use of `GROUP BY` and subqueries in `FROM`, this construct can be avoided, but it is often the most convenient alternative.

7.4. Combining Queries (UNION, INTERSECT, EXCEPT)

The results of two queries can be combined using the set operations union, intersection, and difference. The syntax is

```
query1 UNION [ALL] query2
query1 INTERSECT [ALL] query2
query1 EXCEPT [ALL] query2
```

where *query1* and *query2* are queries that can use any of the features discussed up to this point.

`UNION` effectively appends the result of *query2* to the result of *query1* (although there is no guarantee that this is the order in which the rows are actually returned). Furthermore, it eliminates duplicate rows from its result, in the same way as `DISTINCT`, unless `UNION ALL` is used.

`INTERSECT` returns all rows that are both in the result of *query1* and in the result of *query2*. Duplicate rows are eliminated unless `INTERSECT ALL` is used.

`EXCEPT` returns all rows that are in the result of *query1* but not in the result of *query2*. (This is sometimes called the *difference* between two queries.) Again, duplicates are eliminated unless `EXCEPT ALL` is used.

In order to calculate the union, intersection, or difference of two queries, the two queries must be “union compatible”, which means that they return the same number of columns and the corresponding columns have compatible data types, as described in [Section 10.5](#).

Set operations can be combined, for example

```
query1 UNION query2 EXCEPT query3
```

which is equivalent to

```
(query1 UNION query2) EXCEPT query3
```

As shown here, you can use parentheses to control the order of evaluation. Without parentheses, `UNION` and `EXCEPT` associate left-to-right, but `INTERSECT` binds more tightly than those two operators. Thus

```
query1 UNION query2 INTERSECT query3
```

means

```
query1 UNION (query2 INTERSECT query3)
```

You can also surround an individual *query* with parentheses. This is important if the *query* needs to use any of the clauses discussed in following sections, such as `LIMIT`. Without parentheses, you'll get a syntax error, or else the clause will be understood as applying to the output of the set operation rather than one of its inputs. For example,

```
SELECT a FROM b UNION SELECT x FROM y LIMIT 10
```

is accepted, but it means

```
(SELECT a FROM b UNION SELECT x FROM y) LIMIT 10
```

not

```
SELECT a FROM b UNION (SELECT x FROM y LIMIT 10)
```

7.5. Sorting Rows (`ORDER BY`)

After a query has produced an output table (after the select list has been processed) it can optionally be sorted. If sorting is not chosen, the rows will be returned in an unspecified order. The actual order in that case will depend on the scan and join plan types and the order on disk, but it must not be relied on. A particular output ordering can only be guaranteed if the sort step is explicitly chosen.

The `ORDER BY` clause specifies the sort order:

```
SELECT select_list
   FROM table_expression
   ORDER BY sort_expression1 [ASC | DESC] [NULLS { FIRST | LAST }]
           [, sort_expression2 [ASC | DESC] [NULLS { FIRST | LAST }] ...]
```

The sort expression(s) can be any expression that would be valid in the query's select list. An example is:

```
SELECT a, b FROM table1 ORDER BY a + b, c;
```

When more than one expression is specified, the later values are used to sort rows that are equal according to the earlier values. Each expression can be followed by an optional `ASC` or `DESC` keyword to set the sort direction to ascending or descending. `ASC` order is the default. Ascending order puts smaller values first, where “smaller” is defined in terms of the `<` operator. Similarly, descending order is determined with the `>` operator.¹

The `NULLS FIRST` and `NULLS LAST` options can be used to determine whether nulls appear before or after non-null values in the sort ordering. By default, null values sort as if larger than any non-null value; that is, `NULLS FIRST` is the default for `DESC` order, and `NULLS LAST` otherwise.

Note that the ordering options are considered independently for each sort column. For example `ORDER BY x, y DESC` means `ORDER BY x ASC, y DESC`, which is not the same as `ORDER BY x DESC, y DESC`.

A *sort_expression* can also be the column label or number of an output column, as in:

¹ Actually, Postgres Pro uses the *default B-tree operator class* for the expression's data type to determine the sort ordering for `ASC` and `DESC`. Conventionally, data types will be set up so that the `<` and `>` operators correspond to this sort ordering, but a user-defined data type's designer could choose to do something different.

```
SELECT a + b AS sum, c FROM table1 ORDER BY sum;
SELECT a, max(b) FROM table1 GROUP BY a ORDER BY 1;
```

both of which sort by the first output column. Note that an output column name has to stand alone, that is, it cannot be used in an expression — for example, this is *not* correct:

```
SELECT a + b AS sum, c FROM table1 ORDER BY sum + c;           -- wrong
```

This restriction is made to reduce ambiguity. There is still ambiguity if an `ORDER BY` item is a simple name that could match either an output column name or a column from the table expression. The output column is used in such cases. This would only cause confusion if you use `AS` to rename an output column to match some other table column's name.

`ORDER BY` can be applied to the result of a `UNION`, `INTERSECT`, or `EXCEPT` combination, but in this case it is only permitted to sort by output column names or numbers, not by expressions.

7.6. LIMIT and OFFSET

`LIMIT` and `OFFSET` allow you to retrieve just a portion of the rows that are generated by the rest of the query:

```
SELECT select_list
      FROM table_expression
      [ ORDER BY ... ]
      [ LIMIT { count | ALL } ]
      [ OFFSET start ]
```

If a limit count is given, no more than that many rows will be returned (but possibly fewer, if the query itself yields fewer rows). `LIMIT ALL` is the same as omitting the `LIMIT` clause, as is `LIMIT` with a `NULL` argument.

`OFFSET` says to skip that many rows before beginning to return rows. `OFFSET 0` is the same as omitting the `OFFSET` clause, as is `OFFSET` with a `NULL` argument.

If both `OFFSET` and `LIMIT` appear, then `OFFSET` rows are skipped before starting to count the `LIMIT` rows that are returned.

When using `LIMIT`, it is important to use an `ORDER BY` clause that constrains the result rows into a unique order. Otherwise you will get an unpredictable subset of the query's rows. You might be asking for the tenth through twentieth rows, but tenth through twentieth in what ordering? The ordering is unknown, unless you specified `ORDER BY`.

The query optimizer takes `LIMIT` into account when generating query plans, so you are very likely to get different plans (yielding different row orders) depending on what you give for `LIMIT` and `OFFSET`. Thus, using different `LIMIT/OFFSET` values to select different subsets of a query result *will give inconsistent results* unless you enforce a predictable result ordering with `ORDER BY`. This is not a bug; it is an inherent consequence of the fact that SQL does not promise to deliver the results of a query in any particular order unless `ORDER BY` is used to constrain the order.

The rows skipped by an `OFFSET` clause still have to be computed inside the server; therefore a large `OFFSET` might be inefficient.

7.7. VALUES Lists

`VALUES` provides a way to generate a “constant table” that can be used in a query without having to actually create and populate a table on-disk. The syntax is

```
VALUES ( expression [, ...] ) [, ...]
```

Each parenthesized list of expressions generates a row in the table. The lists must all have the same number of elements (i.e., the number of columns in the table), and corresponding entries in each list must

have compatible data types. The actual data type assigned to each column of the result is determined using the same rules as for `UNION` (see [Section 10.5](#)).

As an example:

```
VALUES (1, 'one'), (2, 'two'), (3, 'three');
```

will return a table of two columns and three rows. It's effectively equivalent to:

```
SELECT 1 AS column1, 'one' AS column2
UNION ALL
SELECT 2, 'two'
UNION ALL
SELECT 3, 'three';
```

By default, Postgres Pro assigns the names `column1`, `column2`, etc. to the columns of a `VALUES` table. The column names are not specified by the SQL standard and different database systems do it differently, so it's usually better to override the default names with a table alias list, like this:

```
=> SELECT * FROM (VALUES (1, 'one'), (2, 'two'), (3, 'three')) AS t (num,letter);
 num | letter
-----+-----
   1 | one
   2 | two
   3 | three
(3 rows)
```

Syntactically, `VALUES` followed by expression lists is treated as equivalent to:

```
SELECT select_list FROM table_expression
```

and can appear anywhere a `SELECT` can. For example, you can use it as part of a `UNION`, or attach a *sort_specification* (`ORDER BY`, `LIMIT`, and/or `OFFSET`) to it. `VALUES` is most commonly used as the data source in an `INSERT` command, and next most commonly as a subquery.

For more information see [VALUES](#).

7.8. WITH Queries (Common Table Expressions)

`WITH` provides a way to write auxiliary statements for use in a larger query. These statements, which are often referred to as Common Table Expressions or CTEs, can be thought of as defining temporary tables that exist just for one query. Each auxiliary statement in a `WITH` clause can be a `SELECT`, `INSERT`, `UPDATE`, `DELETE`, or `MERGE`; and the `WITH` clause itself is attached to a primary statement that can also be a `SELECT`, `INSERT`, `UPDATE`, `DELETE`, or `MERGE`.

7.8.1. SELECT in WITH

The basic value of `SELECT` in `WITH` is to break down complicated queries into simpler parts. An example is:

```
WITH regional_sales AS (
    SELECT region, SUM(amount) AS total_sales
    FROM orders
    GROUP BY region
), top_regions AS (
    SELECT region
    FROM regional_sales
    WHERE total_sales > (SELECT SUM(total_sales)/10 FROM regional_sales)
)
SELECT region,
       product,
       SUM(quantity) AS product_units,
       SUM(amount) AS product_sales
```

```
FROM orders
WHERE region IN (SELECT region FROM top_regions)
GROUP BY region, product;
```

which displays per-product sales totals in only the top sales regions. The `WITH` clause defines two auxiliary statements named `regional_sales` and `top_regions`, where the output of `regional_sales` is used in `top_regions` and the output of `top_regions` is used in the primary `SELECT` query. This example could have been written without `WITH`, but we'd have needed two levels of nested sub-`SELECT`s. It's a bit easier to follow this way.

7.8.2. Recursive Queries

The optional `RECURSIVE` modifier changes `WITH` from a mere syntactic convenience into a feature that accomplishes things not otherwise possible in standard SQL. Using `RECURSIVE`, a `WITH` query can refer to its own output. A very simple example is this query to sum the integers from 1 through 100:

```
WITH RECURSIVE t(n) AS (
    VALUES (1)
    UNION ALL
    SELECT n+1 FROM t WHERE n < 100
)
SELECT sum(n) FROM t;
```

The general form of a recursive `WITH` query is always a *non-recursive term*, then `UNION` (or `UNION ALL`), then a *recursive term*, where only the recursive term can contain a reference to the query's own output. Such a query is executed as follows:

Recursive Query Evaluation

1. Evaluate the non-recursive term. For `UNION` (but not `UNION ALL`), discard duplicate rows. Include all remaining rows in the result of the recursive query, and also place them in a temporary *working table*.
2. So long as the working table is not empty, repeat these steps:
 - a. Evaluate the recursive term, substituting the current contents of the working table for the recursive self-reference. For `UNION` (but not `UNION ALL`), discard duplicate rows and rows that duplicate any previous result row. Include all remaining rows in the result of the recursive query, and also place them in a temporary *intermediate table*.
 - b. Replace the contents of the working table with the contents of the intermediate table, then empty the intermediate table.

Note

While `RECURSIVE` allows queries to be specified recursively, internally such queries are evaluated iteratively.

In the example above, the working table has just a single row in each step, and it takes on the values from 1 through 100 in successive steps. In the 100th step, there is no output because of the `WHERE` clause, and so the query terminates.

Recursive queries are typically used to deal with hierarchical or tree-structured data. A useful example is this query to find all the direct and indirect sub-parts of a product, given only a table that shows immediate inclusions:

```
WITH RECURSIVE included_parts(sub_part, part, quantity) AS (
    SELECT sub_part, part, quantity FROM parts WHERE part = 'our_product'
    UNION ALL
    SELECT p.sub_part, p.part, p.quantity * pr.quantity
    FROM included_parts pr, parts p
    WHERE p.part = pr.sub_part
```

```
)
SELECT sub_part, SUM(quantity) as total_quantity
FROM included_parts
GROUP BY sub_part
```

7.8.2.1. Search Order

When computing a tree traversal using a recursive query, you might want to order the results in either depth-first or breadth-first order. This can be done by computing an ordering column alongside the other data columns and using that to sort the results at the end. Note that this does not actually control in which order the query evaluation visits the rows; that is as always in SQL implementation-dependent. This approach merely provides a convenient way to order the results afterwards.

To create a depth-first order, we compute for each result row an array of rows that we have visited so far. For example, consider the following query that searches a table `tree` using a `link` field:

```
WITH RECURSIVE search_tree(id, link, data) AS (
    SELECT t.id, t.link, t.data
    FROM tree t
    UNION ALL
    SELECT t.id, t.link, t.data
    FROM tree t, search_tree st
    WHERE t.id = st.link
)
SELECT * FROM search_tree;
```

To add depth-first ordering information, you can write this:

```
WITH RECURSIVE search_tree(id, link, data, path) AS (
    SELECT t.id, t.link, t.data, ARRAY[t.id]
    FROM tree t
    UNION ALL
    SELECT t.id, t.link, t.data, path || t.id
    FROM tree t, search_tree st
    WHERE t.id = st.link
)
SELECT * FROM search_tree ORDER BY path;
```

In the general case where more than one field needs to be used to identify a row, use an array of rows. For example, if we needed to track fields `f1` and `f2`:

```
WITH RECURSIVE search_tree(id, link, data, path) AS (
    SELECT t.id, t.link, t.data, ARRAY[ROW(t.f1, t.f2)]
    FROM tree t
    UNION ALL
    SELECT t.id, t.link, t.data, path || ROW(t.f1, t.f2)
    FROM tree t, search_tree st
    WHERE t.id = st.link
)
SELECT * FROM search_tree ORDER BY path;
```

Tip

Omit the `ROW()` syntax in the common case where only one field needs to be tracked. This allows a simple array rather than a composite-type array to be used, gaining efficiency.

To create a breadth-first order, you can add a column that tracks the depth of the search, for example:

```
WITH RECURSIVE search_tree(id, link, data, depth) AS (
    SELECT t.id, t.link, t.data, 0
```

```

FROM tree t
UNION ALL
SELECT t.id, t.link, t.data, depth + 1
FROM tree t, search_tree st
WHERE t.id = st.link
)
SELECT * FROM search_tree ORDER BY depth;

```

To get a stable sort, add data columns as secondary sorting columns.

Tip

The recursive query evaluation algorithm produces its output in breadth-first search order. However, this is an implementation detail and it is perhaps unsound to rely on it. The order of the rows within each level is certainly undefined, so some explicit ordering might be desired in any case.

There is built-in syntax to compute a depth- or breadth-first sort column. For example:

```

WITH RECURSIVE search_tree(id, link, data) AS (
    SELECT t.id, t.link, t.data
    FROM tree t
    UNION ALL
    SELECT t.id, t.link, t.data
    FROM tree t, search_tree st
    WHERE t.id = st.link
) SEARCH DEPTH FIRST BY id SET ordercol
SELECT * FROM search_tree ORDER BY ordercol;

```

```

WITH RECURSIVE search_tree(id, link, data) AS (
    SELECT t.id, t.link, t.data
    FROM tree t
    UNION ALL
    SELECT t.id, t.link, t.data
    FROM tree t, search_tree st
    WHERE t.id = st.link
) SEARCH BREADTH FIRST BY id SET ordercol
SELECT * FROM search_tree ORDER BY ordercol;

```

This syntax is internally expanded to something similar to the above hand-written forms. The `SEARCH` clause specifies whether depth- or breadth first search is wanted, the list of columns to track for sorting, and a column name that will contain the result data that can be used for sorting. That column will implicitly be added to the output rows of the CTE.

7.8.2.2. Cycle Detection

When working with recursive queries it is important to be sure that the recursive part of the query will eventually return no tuples, or else the query will loop indefinitely. Sometimes, using `UNION` instead of `UNION ALL` can accomplish this by discarding rows that duplicate previous output rows. However, often a cycle does not involve output rows that are completely duplicate: it may be necessary to check just one or a few fields to see if the same point has been reached before. The standard method for handling such situations is to compute an array of the already-visited values. For example, consider again the following query that searches a table `graph` using a `link` field:

```

WITH RECURSIVE search_graph(id, link, data, depth) AS (
    SELECT g.id, g.link, g.data, 0
    FROM graph g
    UNION ALL
    SELECT g.id, g.link, g.data, sg.depth + 1
    FROM graph g, search_graph sg

```



```
WHERE g.id = sg.link
)
SELECT * FROM search_graph;
```

This query will loop if the link relationships contain cycles. Because we require a “depth” output, just changing UNION ALL to UNION would not eliminate the looping. Instead we need to recognize whether we have reached the same row again while following a particular path of links. We add two columns `is_cycle` and `path` to the loop-prone query:

```
WITH RECURSIVE search_graph(id, link, data, depth, is_cycle, path) AS (
    SELECT g.id, g.link, g.data, 0,
           false,
           ARRAY[g.id]
    FROM graph g
    UNION ALL
    SELECT g.id, g.link, g.data, sg.depth + 1,
           g.id = ANY(path),
           path || g.id
    FROM graph g, search_graph sg
    WHERE g.id = sg.link AND NOT is_cycle
)
SELECT * FROM search_graph;
```

Aside from preventing cycles, the array value is often useful in its own right as representing the “path” taken to reach any particular row.

In the general case where more than one field needs to be checked to recognize a cycle, use an array of rows. For example, if we needed to compare fields `f1` and `f2`:

```
WITH RECURSIVE search_graph(id, link, data, depth, is_cycle, path) AS (
    SELECT g.id, g.link, g.data, 0,
           false,
           ARRAY[ROW(g.f1, g.f2)]
    FROM graph g
    UNION ALL
    SELECT g.id, g.link, g.data, sg.depth + 1,
           ROW(g.f1, g.f2) = ANY(path),
           path || ROW(g.f1, g.f2)
    FROM graph g, search_graph sg
    WHERE g.id = sg.link AND NOT is_cycle
)
SELECT * FROM search_graph;
```

Tip

Omit the `ROW()` syntax in the common case where only one field needs to be checked to recognize a cycle. This allows a simple array rather than a composite-type array to be used, gaining efficiency.

There is built-in syntax to simplify cycle detection. The above query can also be written like this:

```
WITH RECURSIVE search_graph(id, link, data, depth) AS (
    SELECT g.id, g.link, g.data, 1
    FROM graph g
    UNION ALL
    SELECT g.id, g.link, g.data, sg.depth + 1
    FROM graph g, search_graph sg
    WHERE g.id = sg.link
) CYCLE id SET is_cycle USING path
SELECT * FROM search_graph;
```

and it will be internally rewritten to the above form. The `CYCLE` clause specifies first the list of columns to track for cycle detection, then a column name that will show whether a cycle has been detected, and finally the name of another column that will track the path. The cycle and path columns will implicitly be added to the output rows of the CTE.

Tip

The cycle path column is computed in the same way as the depth-first ordering column show in the previous section. A query can have both a `SEARCH` and a `CYCLE` clause, but a depth-first search specification and a cycle detection specification would create redundant computations, so it's more efficient to just use the `CYCLE` clause and order by the path column. If breadth-first ordering is wanted, then specifying both `SEARCH` and `CYCLE` can be useful.

A helpful trick for testing queries when you are not certain if they might loop is to place a `LIMIT` in the parent query. For example, this query would loop forever without the `LIMIT`:

```
WITH RECURSIVE t(n) AS (
    SELECT 1
    UNION ALL
    SELECT n+1 FROM t
)
SELECT n FROM t LIMIT 100;
```

This works because Postgres Pro's implementation evaluates only as many rows of a `WITH` query as are actually fetched by the parent query. Using this trick in production is not recommended, because other systems might work differently. Also, it usually won't work if you make the outer query sort the recursive query's results or join them to some other table, because in such cases the outer query will usually try to fetch all of the `WITH` query's output anyway.

7.8.3. Common Table Expression Materialization

A useful property of `WITH` queries is that they are normally evaluated only once per execution of the parent query, even if they are referred to more than once by the parent query or sibling `WITH` queries. Thus, expensive calculations that are needed in multiple places can be placed within a `WITH` query to avoid redundant work. Another possible application is to prevent unwanted multiple evaluations of functions with side-effects. However, the other side of this coin is that the optimizer is not able to push restrictions from the parent query down into a multiply-referenced `WITH` query, since that might affect all uses of the `WITH` query's output when it should affect only one. The multiply-referenced `WITH` query will be evaluated as written, without suppression of rows that the parent query might discard afterwards. (But, as mentioned above, evaluation might stop early if the reference(s) to the query demand only a limited number of rows.)

However, if a `WITH` query is non-recursive and side-effect-free (that is, it is a `SELECT` containing no volatile functions) then it can be folded into the parent query, allowing joint optimization of the two query levels. By default, this happens if the parent query references the `WITH` query just once, but not if it references the `WITH` query more than once. You can override that decision by specifying `MATERIALIZED` to force separate calculation of the `WITH` query, or by specifying `NOT MATERIALIZED` to force it to be merged into the parent query. The latter choice risks duplicate computation of the `WITH` query, but it can still give a net savings if each usage of the `WITH` query needs only a small part of the `WITH` query's full output.

A simple example of these rules is

```
WITH w AS (
    SELECT * FROM big_table
)
SELECT * FROM w WHERE key = 123;
```

This `WITH` query will be folded, producing the same execution plan as

```
SELECT * FROM big_table WHERE key = 123;
```

In particular, if there's an index on `key`, it will probably be used to fetch just the rows having `key = 123`. On the other hand, in

```
WITH w AS (
    SELECT * FROM big_table
)
SELECT * FROM w AS w1 JOIN w AS w2 ON w1.key = w2.ref
WHERE w2.key = 123;
```

the `WITH` query will be materialized, producing a temporary copy of `big_table` that is then joined with itself — without benefit of any index. This query will be executed much more efficiently if written as

```
WITH w AS NOT MATERIALIZED (
    SELECT * FROM big_table
)
SELECT * FROM w AS w1 JOIN w AS w2 ON w1.key = w2.ref
WHERE w2.key = 123;
```

so that the parent query's restrictions can be applied directly to scans of `big_table`.

An example where `NOT MATERIALIZED` could be undesirable is

```
WITH w AS (
    SELECT key, very_expensive_function(val) as f FROM some_table
)
SELECT * FROM w AS w1 JOIN w AS w2 ON w1.f = w2.f;
```

Here, materialization of the `WITH` query ensures that `very_expensive_function` is evaluated only once per table row, not twice.

The examples above only show `WITH` being used with `SELECT`, but it can be attached in the same way to `INSERT`, `UPDATE`, `DELETE`, or `MERGE`. In each case it effectively provides temporary table(s) that can be referred to in the main command.

7.8.4. Data-Modifying Statements in `WITH`

You can use data-modifying statements (`INSERT`, `UPDATE`, `DELETE`, or `MERGE`) in `WITH`. This allows you to perform several different operations in the same query. An example is:

```
WITH moved_rows AS (
    DELETE FROM products
    WHERE
        "date" >= '2010-10-01' AND
        "date" < '2010-11-01'
    RETURNING *
)
INSERT INTO products_log
SELECT * FROM moved_rows;
```

This query effectively moves rows from `products` to `products_log`. The `DELETE` in `WITH` deletes the specified rows from `products`, returning their contents by means of its `RETURNING` clause; and then the primary query reads that output and inserts it into `products_log`.

A fine point of the above example is that the `WITH` clause is attached to the `INSERT`, not the sub-`SELECT` within the `INSERT`. This is necessary because data-modifying statements are only allowed in `WITH` clauses that are attached to the top-level statement. However, normal `WITH` visibility rules apply, so it is possible to refer to the `WITH` statement's output from the sub-`SELECT`.

Data-modifying statements in `WITH` usually have `RETURNING` clauses (see [Section 6.4](#)), as shown in the example above. It is the output of the `RETURNING` clause, *not* the target table of the data-modifying statement, that forms the temporary table that can be referred to by the rest of the query. If a data-modifying statement in `WITH` lacks a `RETURNING` clause, then it forms no temporary table and cannot be referred to in the rest of the query. Such a statement will be executed nonetheless. A not-particularly-useful example is:

```
WITH t AS (
    DELETE FROM foo
)
DELETE FROM bar;
```

This example would remove all rows from tables `foo` and `bar`. The number of affected rows reported to the client would only include rows removed from `bar`.

Recursive self-references in data-modifying statements are not allowed. In some cases it is possible to work around this limitation by referring to the output of a recursive `WITH`, for example:

```
WITH RECURSIVE included_parts(sub_part, part) AS (
    SELECT sub_part, part FROM parts WHERE part = 'our_product'
    UNION ALL
    SELECT p.sub_part, p.part
    FROM included_parts pr, parts p
    WHERE p.part = pr.sub_part
)
DELETE FROM parts
WHERE part IN (SELECT part FROM included_parts);
```

This query would remove all direct and indirect subparts of a product.

Data-modifying statements in `WITH` are executed exactly once, and always to completion, independently of whether the primary query reads all (or indeed any) of their output. Notice that this is different from the rule for `SELECT` in `WITH`: as stated in the previous section, execution of a `SELECT` is carried only as far as the primary query demands its output.

The sub-statements in `WITH` are executed concurrently with each other and with the main query. Therefore, when using data-modifying statements in `WITH`, the order in which the specified updates actually happen is unpredictable. All the statements are executed with the same *snapshot* (see [Chapter 13](#)), so they cannot “see” one another’s effects on the target tables. This alleviates the effects of the unpredictability of the actual order of row updates, and means that `RETURNING` data is the only way to communicate changes between different `WITH` sub-statements and the main query. An example of this is that in

```
WITH t AS (
    UPDATE products SET price = price * 1.05
    RETURNING *
)
SELECT * FROM products;
```

the outer `SELECT` would return the original prices before the action of the `UPDATE`, while in

```
WITH t AS (
    UPDATE products SET price = price * 1.05
    RETURNING *
)
SELECT * FROM t;
```

the outer `SELECT` would return the updated data.

Trying to update the same row twice in a single statement is not supported. Only one of the modifications takes place, but it is not easy (and sometimes not possible) to reliably predict which one. This also applies to deleting a row that was already updated in the same statement: only the update is performed. Therefore you should generally avoid trying to modify a single row twice in a single statement. In particular avoid writing `WITH` sub-statements that could affect the same rows changed by the main statement or a sibling sub-statement. The effects of such a statement will not be predictable.

At present, any table used as the target of a data-modifying statement in `WITH` must not have a conditional rule, nor an `ALSO` rule, nor an `INSTEAD` rule that expands to multiple statements.

Chapter 8. Data Types

Postgres Pro has a rich set of native data types available to users. Users can add new types to Postgres Pro using the [CREATE TYPE](#) command.

[Table 8.1](#) shows all the built-in general-purpose data types. Most of the alternative names listed in the “Aliases” column are the names used internally by Postgres Pro for historical reasons. In addition, some internally used or deprecated types are available, but are not listed here.

Table 8.1. Data Types

Name	Aliases	Description
bigint	int8	signed eight-byte integer
bigserial	serial8	autoincrementing eight-byte integer
bit [(n)]		fixed-length bit string
bit varying [(n)]	varbit [(n)]	variable-length bit string
boolean	bool	logical Boolean (true/false)
box		rectangular box on a plane
bytea		binary data (“byte array”)
character [(n)]	char [(n)]	fixed-length character string
character varying [(n)]	varchar [(n)]	variable-length character string
cidr		IPv4 or IPv6 network address
circle		circle on a plane
date		calendar date (year, month, day)
double precision	float8	double precision floating-point number (8 bytes)
inet		IPv4 or IPv6 host address
integer	int, int4	signed four-byte integer
interval [fields] [(p)]		time span
json		textual JSON data
jsonb		binary JSON data, decomposed
line		infinite line on a plane
lseg		line segment on a plane
macaddr		MAC (Media Access Control) address
macaddr8		MAC (Media Access Control) address (EUI-64 format)
money		currency amount
numeric [(p, s)]	decimal [(p, s)]	exact numeric of selectable precision
path		geometric path on a plane
pg_lsn		Postgres Pro Log Sequence Number
pg_snapshot		user-level transaction ID snapshot
point		geometric point on a plane
polygon		closed geometric path on a plane
real	float4	single precision floating-point number (4 bytes)

Name	Aliases	Description
smallint	int2	signed two-byte integer
smallserial	serial2	autoincrementing two-byte integer
serial	serial4	autoincrementing four-byte integer
text		variable-length character string
time [(p)] [without time zone]		time of day (no time zone)
time [(p)] with time zone	timetz	time of day, including time zone
timestamp [(p)] [without time zone]		date and time (no time zone)
timestamp [(p)] with time zone	timestampz	date and time, including time zone
tsquery		text search query
tsvector		text search document
txid_snapshot		user-level transaction ID snapshot (deprecated; see <code>pg_snapshot</code>)
uuid		universally unique identifier
xml		XML data

Compatibility

The following types (or spellings thereof) are specified by SQL: `bigint`, `bit`, `bit varying`, `boolean`, `char`, `character varying`, `character`, `varchar`, `date`, `double precision`, `integer`, `interval`, `numeric`, `decimal`, `real`, `smallint`, `time` (with or without time zone), `timestamp` (with or without time zone), `xml`.

Each data type has an external representation determined by its input and output functions. Many of the built-in types have obvious external formats. However, several types are either unique to Postgres Pro, such as geometric paths, or have several possible formats, such as the date and time types. Some of the input and output functions are not invertible, i.e., the result of an output function might lose accuracy when compared to the original input.

8.1. Numeric Types

Numeric types consist of two-, four-, and eight-byte integers, four- and eight-byte floating-point numbers, and selectable-precision decimals. [Table 8.2](#) lists the available types.

Table 8.2. Numeric Types

Name	Storage Size	Description	Range
smallint	2 bytes	small-range integer	-32768 to +32767
integer	4 bytes	typical choice for integer	-2147483648 to +2147483647
bigint	8 bytes	large-range integer	-9223372036854775808 to +9223372036854775807
decimal	variable	user-specified precision, exact	up to 131072 digits before the decimal point; up to 16383 digits after the decimal point
numeric	variable	user-specified precision, exact	up to 131072 digits before the decimal point; up to

Name	Storage Size	Description	Range
			16383 digits after the decimal point
real	4 bytes	variable-precision, inexact	6 decimal digits precision
double precision	8 bytes	variable-precision, inexact	15 decimal digits precision
smallserial	2 bytes	small autoincrementing integer	1 to 32767
serial	4 bytes	autoincrementing integer	1 to 2147483647
bigserial	8 bytes	large autoincrementing integer	1 to 9223372036854775807

The syntax of constants for the numeric types is described in [Section 4.1.2](#). The numeric types have a full set of corresponding arithmetic operators and functions. Refer to [Chapter 9](#) for more information. The following sections describe the types in detail.

8.1.1. Integer Types

The types `smallint`, `integer`, and `bigint` store whole numbers, that is, numbers without fractional components, of various ranges. Attempts to store values outside of the allowed range will result in an error.

The type `integer` is the common choice, as it offers the best balance between range, storage size, and performance. The `smallint` type is generally only used if disk space is at a premium. The `bigint` type is designed to be used when the range of the `integer` type is insufficient.

SQL only specifies the integer types `integer` (or `int`), `smallint`, and `bigint`. The type names `int2`, `int4`, and `int8` are extensions, which are also used by some other SQL database systems.

8.1.2. Arbitrary Precision Numbers

The type `numeric` can store numbers with a very large number of digits. It is especially recommended for storing monetary amounts and other quantities where exactness is required. Calculations with `numeric` values yield exact results where possible, e.g., addition, subtraction, multiplication. However, calculations on `numeric` values are very slow compared to the integer types, or to the floating-point types described in the next section.

We use the following terms below: The *precision* of a `numeric` is the total count of significant digits in the whole number, that is, the number of digits to both sides of the decimal point. The *scale* of a `numeric` is the count of decimal digits in the fractional part, to the right of the decimal point. So the number 23.5141 has a precision of 6 and a scale of 4. Integers can be considered to have a scale of zero.

Both the maximum precision and the maximum scale of a `numeric` column can be configured. To declare a column of type `numeric` use the syntax:

```
NUMERIC(precision, scale)
```

The precision must be positive, while the scale may be positive or negative (see below). Alternatively:

```
NUMERIC(precision)
```

selects a scale of 0. Specifying:

```
NUMERIC
```

without any precision or scale creates an “unconstrained numeric” column in which numeric values of any length can be stored, up to the implementation limits. A column of this kind will not coerce input values to any particular scale, whereas `numeric` columns with a declared scale will coerce input values to that scale. (The SQL standard requires a default scale of 0, i.e., coercion to integer precision. We find this a bit useless. If you're concerned about portability, always specify the precision and scale explicitly.)

Note

The maximum precision that can be explicitly specified in a `numeric` type declaration is 1000. An unconstrained `numeric` column is subject to the limits described in [Table 8.2](#).

If the scale of a value to be stored is greater than the declared scale of the column, the system will round the value to the specified number of fractional digits. Then, if the number of digits to the left of the decimal point exceeds the declared precision minus the declared scale, an error is raised. For example, a column declared as

```
NUMERIC(3, 1)
```

will round values to 1 decimal place and can store values between -99.9 and 99.9, inclusive.

Beginning in Postgres Pro 15, it is allowed to declare a `numeric` column with a negative scale. Then values will be rounded to the left of the decimal point. The precision still represents the maximum number of non-rounded digits. Thus, a column declared as

```
NUMERIC(2, -3)
```

will round values to the nearest thousand and can store values between -99000 and 99000, inclusive. It is also allowed to declare a scale larger than the declared precision. Such a column can only hold fractional values, and it requires the number of zero digits just to the right of the decimal point to be at least the declared scale minus the declared precision. For example, a column declared as

```
NUMERIC(3, 5)
```

will round values to 5 decimal places and can store values between -0.00999 and 0.00999, inclusive.

Note

Postgres Pro permits the scale in a `numeric` type declaration to be any value in the range -1000 to 1000. However, the SQL standard requires the scale to be in the range 0 to *precision*. Using scales outside that range may not be portable to other database systems.

Numeric values are physically stored without any extra leading or trailing zeroes. Thus, the declared precision and scale of a column are maximums, not fixed allocations. (In this sense the `numeric` type is more akin to `varchar(n)` than to `char(n)`.) The actual storage requirement is two bytes for each group of four decimal digits, plus three to eight bytes overhead.

In addition to ordinary numeric values, the `numeric` type has several special values:

```
Infinity  
-Infinity  
NaN
```

These are adapted from the IEEE 754 standard, and represent “infinity”, “negative infinity”, and “not-a-number”, respectively. When writing these values as constants in an SQL command, you must put quotes around them, for example `UPDATE table SET x = '-Infinity'`. On input, these strings are recognized in a case-insensitive manner. The infinity values can alternatively be spelled `inf` and `-inf`.

The infinity values behave as per mathematical expectations. For example, `Infinity` plus any finite value equals `Infinity`, as does `Infinity` plus `Infinity`; but `Infinity` minus `Infinity` yields `NaN` (not a number), because it has no well-defined interpretation. Note that an infinity can only be stored in an unconstrained `numeric` column, because it notionally exceeds any finite precision limit.

The `NaN` (not a number) value is used to represent undefined calculational results. In general, any operation with a `NaN` input yields another `NaN`. The only exception is when the operation's other inputs are such that the same output would be obtained if the `NaN` were to be replaced by any finite or infinite

numeric value; then, that output value is used for `NaN` too. (An example of this principle is that `NaN` raised to the zero power yields one.)

Note

In most implementations of the “not-a-number” concept, `NaN` is not considered equal to any other numeric value (including `NaN`). In order to allow numeric values to be sorted and used in tree-based indexes, Postgres Pro treats `NaN` values as equal, and greater than all non-`NaN` values.

The types `decimal` and `numeric` are equivalent. Both types are part of the SQL standard.

When rounding values, the `numeric` type rounds ties away from zero, while (on most machines) the `real` and `double precision` types round ties to the nearest even number. For example:

```
SELECT x,
       round(x::numeric) AS num_round,
       round(x::double precision) AS dbl_round
FROM generate_series(-3.5, 3.5, 1) as x;
```

x	num_round	dbl_round
-3.5	-4	-4
-2.5	-3	-2
-1.5	-2	-2
-0.5	-1	-0
0.5	1	0
1.5	2	2
2.5	3	2
3.5	4	4

(8 rows)

8.1.3. Floating-Point Types

The data types `real` and `double precision` are inexact, variable-precision numeric types. On all currently supported platforms, these types are implementations of IEEE Standard 754 for Binary Floating-Point Arithmetic (single and double precision, respectively), to the extent that the underlying processor, operating system, and compiler support it.

Inexact means that some values cannot be converted exactly to the internal format and are stored as approximations, so that storing and retrieving a value might show slight discrepancies. Managing these errors and how they propagate through calculations is the subject of an entire branch of mathematics and computer science and will not be discussed here, except for the following points:

- If you require exact storage and calculations (such as for monetary amounts), use the `numeric` type instead.
- If you want to do complicated calculations with these types for anything important, especially if you rely on certain behavior in boundary cases (infinity, underflow), you should evaluate the implementation carefully.
- Comparing two floating-point values for equality might not always work as expected.

On all currently supported platforms, the `real` type has a range of around $1\text{E-}37$ to $1\text{E+}37$ with a precision of at least 6 decimal digits. The `double precision` type has a range of around $1\text{E-}307$ to $1\text{E+}308$ with a precision of at least 15 digits. Values that are too large or too small will cause an error. Rounding might take place if the precision of an input number is too high. Numbers too close to zero that are not representable as distinct from zero will cause an underflow error.

By default, floating point values are output in text form in their shortest precise decimal representation; the decimal value produced is closer to the true stored binary value than to any other value representable in the same binary precision. (However, the output value is currently never *exactly* midway between two

representable values, in order to avoid a widespread bug where input routines do not properly respect the round-to-nearest-even rule.) This value will use at most 17 significant decimal digits for `float8` values, and at most 9 digits for `float4` values.

Note

This shortest-precise output format is much faster to generate than the historical rounded format.

For compatibility with output generated by older versions of Postgres Pro, and to allow the output precision to be reduced, the [extra_float_digits](#) parameter can be used to select rounded decimal output instead. Setting a value of 0 restores the previous default of rounding the value to 6 (for `float4`) or 15 (for `float8`) significant decimal digits. Setting a negative value reduces the number of digits further; for example -2 would round output to 4 or 13 digits respectively.

Any value of [extra_float_digits](#) greater than 0 selects the shortest-precise format.

Note

Applications that wanted precise values have historically had to set [extra_float_digits](#) to 3 to obtain them. For maximum compatibility between versions, they should continue to do so.

In addition to ordinary numeric values, the floating-point types have several special values:

Infinity
-Infinity
NaN

These represent the IEEE 754 special values “infinity”, “negative infinity”, and “not-a-number”, respectively. When writing these values as constants in an SQL command, you must put quotes around them, for example `UPDATE table SET x = '-Infinity'`. On input, these strings are recognized in a case-insensitive manner. The infinity values can alternatively be spelled `inf` and `-inf`.

Note

IEEE 754 specifies that NaN should not compare equal to any other floating-point value (including NaN). In order to allow floating-point values to be sorted and used in tree-based indexes, Postgres Pro treats NaN values as equal, and greater than all non-NaN values.

Postgres Pro also supports the SQL-standard notations `float` and `float(p)` for specifying inexact numeric types. Here, *p* specifies the minimum acceptable precision in *binary* digits. Postgres Pro accepts `float(1)` to `float(24)` as selecting the real type, while `float(25)` to `float(53)` select double precision. Values of *p* outside the allowed range draw an error. `float` with no precision specified is taken to mean double precision.

8.1.4. Serial Types

Note

This section describes a Postgres Pro-specific way to create an autoincrementing column. Another way is to use the SQL-standard identity column feature, described at [Section 5.3](#).

The data types `smallserial`, `serial` and `bigserial` are not true types, but merely a notational convenience for creating unique identifier columns (similar to the `AUTO_INCREMENT` property supported by some other databases). In the current implementation, specifying:

```
CREATE TABLE tablename (
    colname SERIAL
);
```

is equivalent to specifying:

```
CREATE SEQUENCE tablename_colname_seq AS integer;
CREATE TABLE tablename (
    colname integer NOT NULL DEFAULT nextval('tablename_colname_seq')
);
ALTER SEQUENCE tablename_colname_seq OWNED BY tablename.colname;
```

Thus, we have created an integer column and arranged for its default values to be assigned from a sequence generator. A `NOT NULL` constraint is applied to ensure that a null value cannot be inserted. (In most cases you would also want to attach a `UNIQUE` or `PRIMARY KEY` constraint to prevent duplicate values from being inserted by accident, but this is not automatic.) Lastly, the sequence is marked as “owned by” the column, so that it will be dropped if the column or table is dropped.

Note

Because `smallserial`, `serial` and `bigserial` are implemented using sequences, there may be “holes” or gaps in the sequence of values which appears in the column, even if no rows are ever deleted. A value allocated from the sequence is still “used up” even if a row containing that value is never successfully inserted into the table column. This may happen, for example, if the inserting transaction rolls back. See `nextval()` in [Section 9.17](#) for details.

To insert the next value of the sequence into the `serial` column, specify that the `serial` column should be assigned its default value. This can be done either by excluding the column from the list of columns in the `INSERT` statement, or through the use of the `DEFAULT` key word.

The type names `serial` and `serial4` are equivalent: both create integer columns. The type names `bigserial` and `serial8` work the same way, except that they create a `bigint` column. `bigserial` should be used if you anticipate the use of more than 2^{31} identifiers over the lifetime of the table. The type names `smallserial` and `serial2` also work the same way, except that they create a `smallint` column.

The sequence created for a `serial` column is automatically dropped when the owning column is dropped. You can drop the sequence without dropping the column, but this will force removal of the column default expression.

8.2. Monetary Types

The `money` type stores a currency amount with a fixed fractional precision; see [Table 8.3](#). The fractional precision is determined by the database's `lc_monetary` setting. The range shown in the table assumes there are two fractional digits. Input is accepted in a variety of formats, including integer and floating-point literals, as well as typical currency formatting, such as '\$1,000.00'. Output is generally in the latter form but depends on the locale.

Table 8.3. Monetary Types

Name	Storage Size	Description	Range
money	8 bytes	currency amount	-92233720368547758.08 to +92233720368547758.07

Since the output of this data type is locale-sensitive, it might not work to load `money` data into a database that has a different setting of `lc_monetary`. To avoid problems, before restoring a dump into a new database make sure `lc_monetary` has the same or equivalent value as in the database that was dumped.

Values of the `numeric`, `int`, and `bigint` data types can be cast to `money`. Conversion from the `real` and `double precision` data types can be done by casting to `numeric` first, for example:

```
SELECT '12.34'::float8::numeric::money;
```

However, this is not recommended. Floating point numbers should not be used to handle money due to the potential for rounding errors.

A `money` value can be cast to `numeric` without loss of precision. Conversion to other types could potentially lose precision, and must also be done in two stages:

```
SELECT '52093.89'::money::numeric::float8;
```

Division of a `money` value by an integer value is performed with truncation of the fractional part towards zero. To get a rounded result, divide by a floating-point value, or cast the `money` value to `numeric` before dividing and back to `money` afterwards. (The latter is preferable to avoid risking precision loss.) When a `money` value is divided by another `money` value, the result is `double precision` (i.e., a pure number, not money); the currency units cancel each other out in the division.

8.3. Character Types

Table 8.4. Character Types

Name	Description
<code>character varying(n), varchar(n)</code>	variable-length with limit
<code>character(n), char(n), bpchar(n)</code>	fixed-length, blank-padded
<code>bpchar</code>	variable unlimited length, blank-trimmed
<code>text</code>	variable unlimited length

Table 8.4 shows the general-purpose character types available in Postgres Pro.

SQL defines two primary character types: `character varying(n)` and `character(n)`, where `n` is a positive integer. Both of these types can store strings up to `n` characters (not bytes) in length. An attempt to store a longer string into a column of these types will result in an error, unless the excess characters are all spaces, in which case the string will be truncated to the maximum length. (This somewhat bizarre exception is required by the SQL standard.) However, if one explicitly casts a value to `character varying(n)` or `character(n)`, then an over-length value will be truncated to `n` characters without raising an error. (This too is required by the SQL standard.) If the string to be stored is shorter than the declared length, values of type `character` will be space-padded; values of type `character varying` will simply store the shorter string.

In addition, Postgres Pro provides the `text` type, which stores strings of any length. Although the `text` type is not in the SQL standard, several other SQL database management systems have it as well. `text` is Postgres Pro's native string data type, in that most built-in functions operating on strings are declared to take or return `text` not `character varying`. For many purposes, `character varying` acts as though it were a [domain](#) over `text`.

The type name `varchar` is an alias for `character varying`, while `bpchar` (with length specifier) and `char` are aliases for `character`. The `varchar` and `char` aliases are defined in the SQL standard; `bpchar` is a Postgres Pro extension.

If specified, the length `n` must be greater than zero and cannot exceed 10,485,760. If `character varying` (or `varchar`) is used without length specifier, the type accepts strings of any length. If `bpchar` lacks a length specifier, it also accepts strings of any length, but trailing spaces are semantically insignificant. If `character` (or `char`) lacks a specifier, it is equivalent to `character(1)`.

Values of type `character` are physically padded with spaces to the specified width `n`, and are stored and displayed that way. However, trailing spaces are treated as semantically insignificant and disregarded when comparing two values of type `character`. In collations where whitespace is significant, this behavior can produce unexpected results; for example `SELECT 'a '::CHAR(2) collate "C" < E'a`

`\n'::CHAR(2)` returns true, even though C locale would consider a space to be greater than a newline. Trailing spaces are removed when converting a `character` value to one of the other string types. Note that trailing spaces *are* semantically significant in `character varying` and `text` values, and when using pattern matching, that is `LIKE` and regular expressions.

The characters that can be stored in any of these data types are determined by the database character set, which is selected when the database is created. Regardless of the specific character set, the character with code zero (sometimes called NUL) cannot be stored. For more information refer to [Section 23.3](#).

The storage requirement for a short string (up to 126 bytes) is 1 byte plus the actual string, which includes the space padding in the case of `character`. Longer strings have 4 bytes of overhead instead of 1. Long strings are compressed by the system automatically, so the physical requirement on disk might be less. Very long values are also stored in background tables so that they do not interfere with rapid access to shorter column values. In any case, the longest possible character string that can be stored is about 1 GB. (The maximum value that will be allowed for *n* in the data type declaration is less than that. It wouldn't be useful to change this because with multibyte character encodings the number of characters and bytes can be quite different. If you desire to store long strings with no specific upper limit, use `text` or `character varying` without a length specifier, rather than making up an arbitrary length limit.)

Tip

There is no performance difference among these three types, apart from increased storage space when using the blank-padded type, and a few extra CPU cycles to check the length when storing into a length-constrained column. While `character(n)` has performance advantages in some other database systems, there is no such advantage in Postgres Pro; in fact `character(n)` is usually the slowest of the three because of its additional storage costs. In most situations `text` or `character varying` should be used instead.

Refer to [Section 4.1.2.1](#) for information about the syntax of string literals, and to [Chapter 9](#) for information about available operators and functions.

Example 8.1. Using the Character Types

```
CREATE TABLE test1 (a character(4));
INSERT INTO test1 VALUES ('ok');
SELECT a, char_length(a) FROM test1; -- 1
```

a	char_length
ok	2

```
CREATE TABLE test2 (b varchar(5));
INSERT INTO test2 VALUES ('ok');
INSERT INTO test2 VALUES ('good');
INSERT INTO test2 VALUES ('too long');
ERROR:  value too long for type character varying(5)
INSERT INTO test2 VALUES ('too long'::varchar(5)); -- explicit truncation
SELECT b, char_length(b) FROM test2;
```

b	char_length
ok	2
good	5
too l	5

1 The `char_length` function is discussed in [Section 9.4](#).

There are two other fixed-length character types in Postgres Pro, shown in [Table 8.5](#). These are not intended for general-purpose use, only for use in the internal system catalogs. The `name` type is used to store identifiers. Its length is currently defined as 64 bytes (63 usable characters plus terminator) but should be referenced using the constant `NAMEDATALEN` in C source code. The length is set at compile time (and is therefore adjustable for special uses); the default maximum length might change in a future release. The type `"char"` (note the quotes) is different from `char(1)` in that it only uses one byte of storage, and therefore can store only a single ASCII character. It is used in the system catalogs as a simplistic enumeration type.

Table 8.5. Special Character Types

Name	Storage Size	Description
"char"	1 byte	single-byte internal type
name	64 bytes	internal type for object names

8.4. Binary Data Types

The `bytea` data type allows storage of binary strings; see [Table 8.6](#).

Table 8.6. Binary Data Types

Name	Storage Size	Description
bytea	1 or 4 bytes plus the actual binary string	variable-length binary string

A binary string is a sequence of octets (or bytes). Binary strings are distinguished from character strings in two ways. First, binary strings specifically allow storing octets of value zero and other “non-printable” octets (usually, octets outside the decimal range 32 to 126). Character strings disallow zero octets, and also disallow any other octet values and sequences of octet values that are invalid according to the database's selected character set encoding. Second, operations on binary strings process the actual bytes, whereas the processing of character strings depends on locale settings. In short, binary strings are appropriate for storing data that the programmer thinks of as “raw bytes”, whereas character strings are appropriate for storing text.

The `bytea` type supports two formats for input and output: “hex” format and PostgreSQL's historical “escape” format. Both of these are always accepted on input. The output format depends on the configuration parameter `bytea_output`; the default is hex. (Note that the hex format was introduced in PostgreSQL 9.0; earlier versions and some tools don't understand it.)

The SQL standard defines a different binary string type, called `BLOB` or `BINARY LARGE OBJECT`. The input format is different from `bytea`, but the provided functions and operators are mostly the same.

8.4.1. bytea Hex Format

The “hex” format encodes binary data as 2 hexadecimal digits per byte, most significant nibble first. The entire string is preceded by the sequence `\x` (to distinguish it from the escape format). In some contexts, the initial backslash may need to be escaped by doubling it (see [Section 4.1.2.1](#)). For input, the hexadecimal digits can be either upper or lower case, and whitespace is permitted between digit pairs (but not within a digit pair nor in the starting `\x` sequence). The hex format is compatible with a wide range of external applications and protocols, and it tends to be faster to convert than the escape format, so its use is preferred.

Example:

```
SET bytea_output = 'hex';

SELECT '\xDEADBEEF'::bytea;
      bytea
-----
```

```
\xdeadbeef
```

8.4.2. bytea Escape Format

The “escape” format is the traditional Postgres Pro format for the `bytea` type. It takes the approach of representing a binary string as a sequence of ASCII characters, while converting those bytes that cannot be represented as an ASCII character into special escape sequences. If, from the point of view of the application, representing bytes as characters makes sense, then this representation can be convenient. But in practice it is usually confusing because it fuzzes up the distinction between binary strings and character strings, and also the particular escape mechanism that was chosen is somewhat unwieldy. Therefore, this format should probably be avoided for most new applications.

When entering `bytea` values in escape format, octets of certain values *must* be escaped, while all octet values *can* be escaped. In general, to escape an octet, convert it into its three-digit octal value and precede it by a backslash. Backslash itself (octet decimal value 92) can alternatively be represented by double backslashes. [Table 8.7](#) shows the characters that must be escaped, and gives the alternative escape sequences where applicable.

Table 8.7. bytea Literal Escaped Octets

Decimal Octet Value	Description	Escaped Input Representation	Example	Hex Representation
0	zero octet	'\000'	'\000'::bytea	\x00
39	single quote	''' or '\047'	'''::bytea	\x27
92	backslash	'\\' or '\134'	'\\'::bytea	\x5c
0 to 31 and 127 to 255	“non-printable” octets	'\xxx' (octal value)	'\001'::bytea	\x01

The requirement to escape *non-printable* octets varies depending on locale settings. In some instances you can get away with leaving them unescaped.

The reason that single quotes must be doubled, as shown in [Table 8.7](#), is that this is true for any string literal in an SQL command. The generic string-literal parser consumes the outermost single quotes and reduces any pair of single quotes to one data character. What the `bytea` input function sees is just one single quote, which it treats as a plain data character. However, the `bytea` input function treats backslashes as special, and the other behaviors shown in [Table 8.7](#) are implemented by that function.

In some contexts, backslashes must be doubled compared to what is shown above, because the generic string-literal parser will also reduce pairs of backslashes to one data character; see [Section 4.1.2.1](#).

`Bytea` octets are output in hex format by default. If you change `bytea_output` to `escape`, “non-printable” octets are converted to their equivalent three-digit octal value and preceded by one backslash. Most “printable” octets are output by their standard representation in the client character set, e.g.:

```
SET bytea_output = 'escape';

SELECT 'abc \153\154\155 \052\251\124'::bytea;
      bytea
-----
abc klm *\251T
```

The octet with decimal value 92 (backslash) is doubled in the output. Details are in [Table 8.8](#).

Table 8.8. bytea Output Escaped Octets

Decimal Octet Value	Description	Escaped Output Representation	Example	Output Result
92	backslash	\\	'\134'::bytea	\\

Decimal Octet Value	Description	Escaped Output Representation	Example	Output Result
0 to 31 and 127 to 255	“non-printable” octets	\xxx (octal value)	'\001':::bytea	\001
32 to 126	“printable” octets	client character set representation	'\176':::bytea	~

Depending on the front end to Postgres Pro you use, you might have additional work to do in terms of escaping and unescaping `bytea` strings. For example, you might also have to escape line feeds and carriage returns if your interface automatically translates these.

8.5. Date/Time Types

Postgres Pro supports the full set of SQL date and time types, shown in [Table 8.9](#). The operations available on these data types are described in [Section 9.9](#). Dates are counted according to the Gregorian calendar, even in years before that calendar was introduced (see [Section B.6](#) for more information).

Table 8.9. Date/Time Types

Name	Storage Size	Description	Low Value	High Value	Resolution
<code>timestamp</code> [(<i>p</i>)] [without time zone]	8 bytes	both date and time (no time zone)	4713 BC	294276 AD	1 microsecond
<code>timestamp</code> [(<i>p</i>)] with time zone	8 bytes	both date and time, with time zone	4713 BC	294276 AD	1 microsecond
<code>date</code>	4 bytes	date (no time of day)	4713 BC	5874897 AD	1 day
<code>time</code> [(<i>p</i>)] [without time zone]	8 bytes	time of day (no date)	00:00:00	24:00:00	1 microsecond
<code>time</code> [(<i>p</i>)] with time zone	12 bytes	time of day (no date), with time zone	00:00:00+1559	24:00:00-1559	1 microsecond
<code>interval</code> [<i>fields</i>] [(<i>p</i>)]	16 bytes	time interval	-178000000 years	178000000 years	1 microsecond

Note

The SQL standard requires that writing just `timestamp` be equivalent to `timestamp without time zone`, and Postgres Pro honors that behavior. `timestamptz` is accepted as an abbreviation for `timestamp with time zone`; this is a Postgres Pro extension.

`time`, `timestamp`, and `interval` accept an optional precision value *p* which specifies the number of fractional digits retained in the seconds field. By default, there is no explicit bound on precision. The allowed range of *p* is from 0 to 6.

The `interval` type has an additional option, which is to restrict the set of stored fields by writing one of these phrases:

YEAR

MONTH
 DAY
 HOUR
 MINUTE
 SECOND
 YEAR TO MONTH
 DAY TO HOUR
 DAY TO MINUTE
 DAY TO SECOND
 HOUR TO MINUTE
 HOUR TO SECOND
 MINUTE TO SECOND

Note that if both *fields* and *p* are specified, the *fields* must include `SECOND`, since the precision applies only to the seconds.

The type `time` with `time zone` is defined by the SQL standard, but the definition exhibits properties which lead to questionable usefulness. In most cases, a combination of `date`, `time`, `timestamp without time zone`, and `timestamp with time zone` should provide a complete range of date/time functionality required by any application.

8.5.1. Date/Time Input

Date and time input is accepted in almost any reasonable format, including ISO 8601, SQL-compatible, traditional POSTGRES, and others. For some formats, ordering of day, month, and year in date input is ambiguous and there is support for specifying the expected ordering of these fields. Set the [DateStyle](#) parameter to `MDY` to select month-day-year interpretation, `DMY` to select day-month-year interpretation, or `YMD` to select year-month-day interpretation.

Postgres Pro is more flexible in handling date/time input than the SQL standard requires. See [Appendix B](#) for the exact parsing rules of date/time input and for the recognized text fields including months, days of the week, and time zones.

Remember that any date or time literal input needs to be enclosed in single quotes, like text strings. Refer to [Section 4.1.2.7](#) for more information. SQL requires the following syntax

```
type [ (p) ] 'value'
```

where *p* is an optional precision specification giving the number of fractional digits in the seconds field. Precision can be specified for `time`, `timestamp`, and `interval` types, and can range from 0 to 6. If no precision is specified in a constant specification, it defaults to the precision of the literal value (but not more than 6 digits).

8.5.1.1. Dates

[Table 8.10](#) shows some possible inputs for the `date` type.

Table 8.10. Date Input

Example	Description
1999-01-08	ISO 8601; January 8 in any mode (recommended format)
January 8, 1999	unambiguous in any <code>datestyle</code> input mode
1/8/1999	January 8 in <code>MDY</code> mode; August 1 in <code>DMY</code> mode
1/18/1999	January 18 in <code>MDY</code> mode; rejected in other modes
01/02/03	January 2, 2003 in <code>MDY</code> mode; February 1, 2003 in <code>DMY</code> mode; February 3, 2001 in <code>YMD</code> mode
1999-Jan-08	January 8 in any mode
Jan-08-1999	January 8 in any mode

Example	Description
08-Jan-1999	January 8 in any mode
99-Jan-08	January 8 in YMD mode, else error
08-Jan-99	January 8, except error in YMD mode
Jan-08-99	January 8, except error in YMD mode
19990108	ISO 8601; January 8, 1999 in any mode
990108	ISO 8601; January 8, 1999 in any mode
1999.008	year and day of year
J2451187	Julian date
January 8, 99 BC	year 99 BC

8.5.1.2. Times

The time-of-day types are `time [(p)]` without time zone and `time [(p)]` with time zone. `time` alone is equivalent to `time without time zone`.

Valid input for these types consists of a time of day followed by an optional time zone. (See [Table 8.11](#) and [Table 8.12](#).) If a time zone is specified in the input for `time without time zone`, it is silently ignored. You can also specify a date but it will be ignored, except when you use a time zone name that involves a daylight-savings rule, such as `America/New_York`. In this case specifying the date is required in order to determine whether standard or daylight-savings time applies. The appropriate time zone offset is recorded in the `time with time zone` value and is output as stored; it is not adjusted to the active time zone.

Table 8.11. Time Input

Example	Description
04:05:06.789	ISO 8601
04:05:06	ISO 8601
04:05	ISO 8601
040506	ISO 8601
04:05 AM	same as 04:05; AM does not affect value
04:05 PM	same as 16:05; input hour must be <= 12
04:05:06.789-8	ISO 8601, with time zone as UTC offset
04:05:06-08:00	ISO 8601, with time zone as UTC offset
04:05-08:00	ISO 8601, with time zone as UTC offset
040506-08	ISO 8601, with time zone as UTC offset
040506+0730	ISO 8601, with fractional-hour time zone as UTC offset
040506+07:30:00	UTC offset specified to seconds (not allowed in ISO 8601)
04:05:06 PST	time zone specified by abbreviation
2003-04-12 04:05:06 America/New_York	time zone specified by full name

Table 8.12. Time Zone Input

Example	Description
PST	Abbreviation (for Pacific Standard Time)

Example	Description
America/New_York	Full time zone name
PST8PDT	POSIX-style time zone specification
-8:00:00	UTC offset for PST
-8:00	UTC offset for PST (ISO 8601 extended format)
-800	UTC offset for PST (ISO 8601 basic format)
-8	UTC offset for PST (ISO 8601 basic format)
zulu	Military abbreviation for UTC
z	Short form of <code>zulu</code> (also in ISO 8601)

Refer to [Section 8.5.3](#) for more information on how to specify time zones.

8.5.1.3. Time Stamps

Valid input for the time stamp types consists of the concatenation of a date and a time, followed by an optional time zone, followed by an optional `AD` or `BC`. (Alternatively, `AD/BC` can appear before the time zone, but this is not the preferred ordering.) Thus:

```
1999-01-08 04:05:06
```

and:

```
1999-01-08 04:05:06 -8:00
```

are valid values, which follow the ISO 8601 standard. In addition, the common format:

```
January 8 04:05:06 1999 PST
```

is supported.

The SQL standard differentiates `timestamp without time zone` and `timestamp with time zone` literals by the presence of a “+” or “-” symbol and time zone offset after the time. Hence, according to the standard,

```
TIMESTAMP '2004-10-19 10:23:54'
```

is a `timestamp without time zone`, while

```
TIMESTAMP '2004-10-19 10:23:54+02'
```

is a `timestamp with time zone`. Postgres Pro never examines the content of a literal string before determining its type, and therefore will treat both of the above as `timestamp without time zone`. To ensure that a literal is treated as `timestamp with time zone`, give it the correct explicit type:

```
TIMESTAMP WITH TIME ZONE '2004-10-19 10:23:54+02'
```

In a value that has been determined to be `timestamp without time zone`, Postgres Pro will silently ignore any time zone indication. That is, the resulting value is derived from the date/time fields in the input string, and is not adjusted for time zone.

For `timestamp with time zone` values, an input string that includes an explicit time zone will be converted to UTC (*Universal Coordinated Time*) using the appropriate offset for that time zone. If no time zone is stated in the input string, then it is assumed to be in the time zone indicated by the system's [TimeZone](#) parameter, and is converted to UTC using the offset for the `timezone` zone. In either case, the value is stored internally as UTC, and the originally stated or assumed time zone is not retained.

When a `timestamp with time zone` value is output, it is always converted from UTC to the current `timezone` zone, and displayed as local time in that zone. To see the time in another time zone, either change `timezone` or use the `AT TIME ZONE` construct (see [Section 9.9.4](#)).

Conversions between `timestamp without time zone` and `timestamp with time zone` normally assume that the `timestamp without time zone` value should be taken or given as `timezone local time`. A different time zone can be specified for the conversion using `AT TIME ZONE`.

8.5.1.4. Special Values

Postgres Pro supports several special date/time input values for convenience, as shown in [Table 8.13](#). The values `infinity` and `-infinity` are specially represented inside the system and will be displayed unchanged; but the others are simply notational shorthands that will be converted to ordinary date/time values when read. (In particular, `now` and related strings are converted to a specific time value as soon as they are read.) All of these values need to be enclosed in single quotes when used as constants in SQL commands.

Table 8.13. Special Date/Time Inputs

Input String	Valid Types	Description
<code>epoch</code>	<code>date</code> , <code>timestamp</code>	1970-01-01 00:00:00+00 (Unix system time zero)
<code>infinity</code>	<code>date</code> , <code>timestamp</code> , <code>interval</code>	later than all other time stamps
<code>-infinity</code>	<code>date</code> , <code>timestamp</code> , <code>interval</code>	earlier than all other time stamps
<code>now</code>	<code>date</code> , <code>time</code> , <code>timestamp</code>	current transaction's start time
<code>today</code>	<code>date</code> , <code>timestamp</code>	midnight (00:00) today
<code>tomorrow</code>	<code>date</code> , <code>timestamp</code>	midnight (00:00) tomorrow
<code>yesterday</code>	<code>date</code> , <code>timestamp</code>	midnight (00:00) yesterday
<code>allballs</code>	<code>time</code>	00:00:00.00 UTC

The following SQL-compatible functions can also be used to obtain the current time value for the corresponding data type: `CURRENT_DATE`, `CURRENT_TIME`, `CURRENT_TIMESTAMP`, `LOCALTIME`, `LOCALTIMESTAMP`. (See [Section 9.9.5](#).) Note that these are SQL functions and are *not* recognized in data input strings.

Caution

While the input strings `now`, `today`, `tomorrow`, and `yesterday` are fine to use in interactive SQL commands, they can have surprising behavior when the command is saved to be executed later, for example in prepared statements, views, and function definitions. The string can be converted to a specific time value that continues to be used long after it becomes stale. Use one of the SQL functions instead in such contexts. For example, `CURRENT_DATE + 1` is safer than `'tomorrow'::date`.

8.5.2. Date/Time Output

The output format of the date/time types can be set to one of the four styles ISO 8601, SQL (Ingres), traditional POSTGRES (Unix date format), or German. The default is the ISO format. (The SQL standard requires the use of the ISO 8601 format. The name of the “SQL” output format is a historical accident.) [Table 8.14](#) shows examples of each output style. The output of the `date` and `time` types is generally only the date or time part in accordance with the given examples. However, the POSTGRES style outputs date-only values in ISO format.

Table 8.14. Date/Time Output Styles

Style Specification	Description	Example
ISO	ISO 8601, SQL standard	1997-12-17 07:37:16-08
SQL	traditional style	12/17/1997 07:37:16.00 PST
Postgres	original style	Wed Dec 17 07:37:16 1997 PST

Style Specification	Description	Example
German	regional style	17.12.1997 07:37:16.00 PST

Note

ISO 8601 specifies the use of uppercase letter `T` to separate the date and time. Postgres Pro accepts that format on input, but on output it uses a space rather than `T`, as shown above. This is for readability and for consistency with [RFC 3339](#) as well as some other database systems.

In the SQL and POSTGRES styles, day appears before month if DMY field ordering has been specified, otherwise month appears before day. (See [Section 8.5.1](#) for how this setting also affects interpretation of input values.) [Table 8.15](#) shows examples.

Table 8.15. Date Order Conventions

datestyle Setting	Input Ordering	Example Output
SQL, DMY	<i>day/month/year</i>	17/12/1997 15:37:16.00 CET
SQL, MDY	<i>month/day/year</i>	12/17/1997 07:37:16.00 PST
Postgres, DMY	<i>day/month/year</i>	Wed 17 Dec 07:37:16 1997 PST

In the ISO style, the time zone is always shown as a signed numeric offset from UTC, with positive sign used for zones east of Greenwich. The offset will be shown as *hh* (hours only) if it is an integral number of hours, else as *hh:mm* if it is an integral number of minutes, else as *hh:mm:ss*. (The third case is not possible with any modern time zone standard, but it can appear when working with timestamps that predate the adoption of standardized time zones.) In the other date styles, the time zone is shown as an alphabetic abbreviation if one is in common use in the current zone. Otherwise it appears as a signed numeric offset in ISO 8601 basic format (*hh* or *hhmm*).

The date/time style can be selected by the user using the `SET datestyle` command, the [DateStyle](#) parameter in the `postgresql.conf` configuration file, or the `PGDATESTYLE` environment variable on the server or client.

The formatting function `to_char` (see [Section 9.8](#)) is also available as a more flexible way to format date/time output.

8.5.3. Time Zones

Time zones, and time-zone conventions, are influenced by political decisions, not just earth geometry. Time zones around the world became somewhat standardized during the 1900s, but continue to be prone to arbitrary changes, particularly with respect to daylight-savings rules. Postgres Pro uses the widely-used IANA (Olson) time zone database for information about historical time zone rules. For times in the future, the assumption is that the latest known rules for a given time zone will continue to be observed indefinitely far into the future.

Postgres Pro endeavors to be compatible with the SQL standard definitions for typical usage. However, the SQL standard has an odd mix of date and time types and capabilities. Two obvious problems are:

- Although the `date` type cannot have an associated time zone, the `time` type can. Time zones in the real world have little meaning unless associated with a date as well as a time, since the offset can vary through the year with daylight-saving time boundaries.
- The default time zone is specified as a constant numeric offset from UTC. It is therefore impossible to adapt to daylight-saving time when doing date/time arithmetic across DST boundaries.

To address these difficulties, we recommend using date/time types that contain both date and time when using time zones. We do *not* recommend using the type `time with time zone` (though it is supported by Postgres Pro for legacy applications and for compliance with the SQL standard). Postgres Pro assumes your local time zone for any type containing only date or time.

All timezone-aware dates and times are stored internally in UTC. They are converted to local time in the zone specified by the [TimeZone](#) configuration parameter before being displayed to the client.

Postgres Pro allows you to specify time zones in three different forms:

- A full time zone name, for example `America/New_York`. The recognized time zone names are listed in the `pg_timezone_names` view (see [Section 55.36](#)). Postgres Pro uses the widely-used IANA time zone data for this purpose, so the same time zone names are also recognized by other software.
- A time zone abbreviation, for example `PST`. Such a specification merely defines a particular offset from UTC, in contrast to full time zone names which can imply a set of daylight savings transition rules as well. The recognized abbreviations are listed in the `pg_timezone_abbrevs` view (see [Section 55.35](#)). You cannot set the configuration parameters [TimeZone](#) or [log_timezone](#) to a time zone abbreviation, but you can use abbreviations in date/time input values and with the `AT TIME ZONE` operator.
- In addition to the timezone names and abbreviations, Postgres Pro will accept POSIX-style time zone specifications, as described in [Section B.5](#). This option is not normally preferable to using a named time zone, but it may be necessary if no suitable IANA time zone entry is available.

In short, this is the difference between abbreviations and full names: abbreviations represent a specific offset from UTC, whereas many of the full names imply a local daylight-savings time rule, and so have two possible UTC offsets. As an example, `2014-06-04 12:00 America/New_York` represents noon local time in New York, which for this particular date was Eastern Daylight Time (UTC-4). So `2014-06-04 12:00 EDT` specifies that same time instant. But `2014-06-04 12:00 EST` specifies noon Eastern Standard Time (UTC-5), regardless of whether daylight savings was nominally in effect on that date.

To complicate matters, some jurisdictions have used the same timezone abbreviation to mean different UTC offsets at different times; for example, in Moscow `MSK` has meant UTC+3 in some years and UTC+4 in others. Postgres Pro interprets such abbreviations according to whatever they meant (or had most recently meant) on the specified date; but, as with the `EST` example above, this is not necessarily the same as local civil time on that date.

In all cases, timezone names and abbreviations are recognized case-insensitively. (This is a change from PostgreSQL versions prior to 8.2, which were case-sensitive in some contexts but not others.)

Neither timezone names nor abbreviations are hard-wired into the server; they are obtained from configuration files stored under `.../share/timezone/` and `.../share/timezonesets/` of the installation directory (see [Section B.4](#)).

The [TimeZone](#) configuration parameter can be set in the file `postgresql.conf`, or in any of the other standard ways described in [Chapter 19](#). There are also some special ways to set it:

- The SQL command `SET TIME ZONE` sets the time zone for the session. This is an alternative spelling of `SET TIMEZONE TO` with a more SQL-spec-compatible syntax.
- The `PGTZ` environment variable is used by libpq clients to send a `SET TIME ZONE` command to the server upon connection.

8.5.4. Interval Input

interval values can be written using the following verbose syntax:

```
[@] quantity unit [quantity unit...] [direction]
```

where *quantity* is a number (possibly signed); *unit* is `microsecond`, `millisecond`, `second`, `minute`, `hour`, `day`, `week`, `month`, `year`, `decade`, `century`, `millennium`, or abbreviations or plurals of these units; *direction* can be `ago` or empty. The at sign (`@`) is optional noise. The amounts of the different units are implicitly added with appropriate sign accounting. `ago` negates all the fields. This syntax is also used for interval output, if [IntervalStyle](#) is set to `postgres_verbose`.

Quantities of days, hours, minutes, and seconds can be specified without explicit unit markings. For example, `'1 12:59:10'` is read the same as `'1 day 12 hours 59 min 10 sec'`. Also, a combination of

years and months can be specified with a dash; for example '200-10' is read the same as '200 years 10 months'. (These shorter forms are in fact the only ones allowed by the SQL standard, and are used for output when `IntervalStyle` is set to `sql_standard`.)

Interval values can also be written as ISO 8601 time intervals, using either the “format with designators” of the standard's section 4.4.3.2 or the “alternative format” of section 4.4.3.3. The format with designators looks like this:

```
P quantity unit [ quantity unit ...] [ T [ quantity unit ...]]
```

The string must start with a `P`, and may include a `T` that introduces the time-of-day units. The available unit abbreviations are given in [Table 8.16](#). Units may be omitted, and may be specified in any order, but units smaller than a day must appear after `T`. In particular, the meaning of `M` depends on whether it is before or after `T`.

Table 8.16. ISO 8601 Interval Unit Abbreviations

Abbreviation	Meaning
Y	Years
M	Months (in the date part)
W	Weeks
D	Days
H	Hours
M	Minutes (in the time part)
S	Seconds

In the alternative format:

```
P [ years-months-days ] [ T hours:minutes:seconds ]
```

the string must begin with `P`, and a `T` separates the date and time parts of the interval. The values are given as numbers similar to ISO 8601 dates.

When writing an interval constant with a *fields* specification, or when assigning a string to an interval column that was defined with a *fields* specification, the interpretation of unmarked quantities depends on the *fields*. For example `INTERVAL '1' YEAR` is read as 1 year, whereas `INTERVAL '1'` means 1 second. Also, field values “to the right” of the least significant field allowed by the *fields* specification are silently discarded. For example, writing `INTERVAL '1 day 2:03:04' HOUR TO MINUTE` results in dropping the seconds field, but not the day field.

According to the SQL standard all fields of an interval value must have the same sign, so a leading negative sign applies to all fields; for example the negative sign in the interval literal `'-1 2:03:04'` applies to both the days and hour/minute/second parts. Postgres Pro allows the fields to have different signs, and traditionally treats each field in the textual representation as independently signed, so that the hour/minute/second part is considered positive in this example. If `IntervalStyle` is set to `sql_standard` then a leading sign is considered to apply to all fields (but only if no additional signs appear). Otherwise the traditional Postgres Pro interpretation is used. To avoid ambiguity, it's recommended to attach an explicit sign to each field if any field is negative.

Internally, interval values are stored as three integral fields: months, days, and microseconds. These fields are kept separate because the number of days in a month varies, while a day can have 23 or 25 hours if a daylight savings time transition is involved. An interval input string that uses other units is normalized into this format, and then reconstructed in a standardized way for output, for example:

```
SELECT '2 years 15 months 100 weeks 99 hours 123456789 milliseconds'::interval;
          interval
-----
3 years 3 mons 700 days 133:17:36.789
```


Here weeks, which are understood as “7 days”, have been kept separate, while the smaller and larger time units were combined and normalized.

Input field values can have fractional parts, for example '1.5 weeks' or '01:02:03.45'. However, because `interval` internally stores only integral fields, fractional values must be converted into smaller units. Fractional parts of units greater than months are rounded to be an integer number of months, e.g. '1.5 years' becomes '1 year 6 mons'. Fractional parts of weeks and days are computed to be an integer number of days and microseconds, assuming 30 days per month and 24 hours per day, e.g., '1.75 months' becomes 1 mon 22 days 12:00:00. Only seconds will ever be shown as fractional on output.

Table 8.17 shows some examples of valid `interval` input.

Table 8.17. Interval Input

Example	Description
1-2	SQL standard format: 1 year 2 months
3 4:05:06	SQL standard format: 3 days 4 hours 5 minutes 6 seconds
1 year 2 months 3 days 4 hours 5 minutes 6 seconds	Traditional Postgres format: 1 year 2 months 3 days 4 hours 5 minutes 6 seconds
P1Y2M3DT4H5M6S	ISO 8601 “format with designators”: same meaning as above
P0001-02-03T04:05:06	ISO 8601 “alternative format”: same meaning as above

8.5.5. Interval Output

As previously explained, Postgres Pro stores `interval` values as months, days, and microseconds. For output, the months field is converted to years and months by dividing by 12. The days field is shown as-is. The microseconds field is converted to hours, minutes, seconds, and fractional seconds. Thus months, minutes, and seconds will never be shown as exceeding the ranges 0-11, 0-59, and 0-59 respectively, while the displayed years, days, and hours fields can be quite large. (The `justify_days` and `justify_hours` functions can be used if it is desirable to transpose large days or hours values into the next higher field.)

The output format of the interval type can be set to one of the four styles `sql_standard`, `postgres`, `postgres_verbose`, or `iso_8601`, using the command `SET intervalstyle`. The default is the `postgres` format. Table 8.18 shows examples of each output style.

The `sql_standard` style produces output that conforms to the SQL standard's specification for interval literal strings, if the interval value meets the standard's restrictions (either year-month only or day-time only, with no mixing of positive and negative components). Otherwise the output looks like a standard year-month literal string followed by a day-time literal string, with explicit signs added to disambiguate mixed-sign intervals.

The output of the `postgres` style matches the output of PostgreSQL releases prior to 8.4 when the `DateStyle` parameter was set to `ISO`.

The output of the `postgres_verbose` style matches the output of PostgreSQL releases prior to 8.4 when the `DateStyle` parameter was set to `non-ISO` output.

The output of the `iso_8601` style matches the “format with designators” described in section 4.4.3.2 of the ISO 8601 standard.

Table 8.18. Interval Output Style Examples

Style Specification	Year-Month Interval	Day-Time Interval	Mixed Interval
<code>sql_standard</code>	1-2	3 4:05:06	-1-2 +3 -4:05:06

Style Specification	Year-Month Interval	Day-Time Interval	Mixed Interval
postgres	1 year 2 mons	3 days 04:05:06	-1 year -2 mons +3 days -04:05:06
postgres_verbose	@ 1 year 2 mons	@ 3 days 4 hours 5 mins 6 secs	@ 1 year 2 mons -3 days 4 hours 5 mins 6 secs ago
iso_8601	P1Y2M	P3DT4H5M6S	P-1Y-2M3DT-4H-5M-6S

8.6. Boolean Type

Postgres Pro provides the standard SQL type `boolean`; see [Table 8.19](#). The `boolean` type can have several states: “true”, “false”, and a third state, “unknown”, which is represented by the SQL null value.

Table 8.19. Boolean Data Type

Name	Storage Size	Description
<code>boolean</code>	1 byte	state of true or false

Boolean constants can be represented in SQL queries by the SQL key words `TRUE`, `FALSE`, and `NULL`.

The datatype input function for type `boolean` accepts these string representations for the “true” state:

```
true
yes
on
1
```

and these representations for the “false” state:

```
false
no
off
0
```

Unique prefixes of these strings are also accepted, for example `t` or `n`. Leading or trailing whitespace is ignored, and case does not matter.

The datatype output function for type `boolean` always emits either `t` or `f`, as shown in [Example 8.2](#).

Example 8.2. Using the `boolean` Type

```
CREATE TABLE test1 (a boolean, b text);
INSERT INTO test1 VALUES (TRUE, 'sic est');
INSERT INTO test1 VALUES (FALSE, 'non est');
SELECT * FROM test1;
 a |      b
---+-----
 t | sic est
 f | non est

SELECT * FROM test1 WHERE a;
 a |      b
---+-----
 t | sic est
```

The key words `TRUE` and `FALSE` are the preferred (SQL-compliant) method for writing Boolean constants in SQL queries. But you can also use the string representations by following the generic string-literal constant syntax described in [Section 4.1.2.7](#), for example `'yes'::boolean`.

Note that the parser automatically understands that `TRUE` and `FALSE` are of type `boolean`, but this is not so for `NULL` because that can have any type. So in some contexts you might have to cast `NULL` to `boolean` explicitly, for example `NULL::boolean`. Conversely, the cast can be omitted from a string-literal `Boolean` value in contexts where the parser can deduce that the literal must be of type `boolean`.

8.7. Enumerated Types

Enumerated (enum) types are data types that comprise a static, ordered set of values. They are equivalent to the `enum` types supported in a number of programming languages. An example of an enum type might be the days of the week, or a set of status values for a piece of data.

8.7.1. Declaration of Enumerated Types

Enum types are created using the `CREATE TYPE` command, for example:

```
CREATE TYPE mood AS ENUM ('sad', 'ok', 'happy');
```

Once created, the enum type can be used in table and function definitions much like any other type:

```
CREATE TYPE mood AS ENUM ('sad', 'ok', 'happy');
CREATE TABLE person (
    name text,
    current_mood mood
);
INSERT INTO person VALUES ('Moe', 'happy');
SELECT * FROM person WHERE current_mood = 'happy';
 name | current_mood
-----+-----
 Moe  | happy
(1 row)
```

8.7.2. Ordering

The ordering of the values in an enum type is the order in which the values were listed when the type was created. All standard comparison operators and related aggregate functions are supported for enums. For example:

```
INSERT INTO person VALUES ('Larry', 'sad');
INSERT INTO person VALUES ('Curly', 'ok');
SELECT * FROM person WHERE current_mood > 'sad';
 name | current_mood
-----+-----
 Moe  | happy
 Curly | ok
(2 rows)

SELECT * FROM person WHERE current_mood > 'sad' ORDER BY current_mood;
 name | current_mood
-----+-----
 Curly | ok
 Moe  | happy
(2 rows)

SELECT name
FROM person
WHERE current_mood = (SELECT MIN(current_mood) FROM person);
 name
-----
 Larry
(1 row)
```

8.7.3. Type Safety

Each enumerated data type is separate and cannot be compared with other enumerated types. See this example:

```
CREATE TYPE happiness AS ENUM ('happy', 'very happy', 'ecstatic');
CREATE TABLE holidays (
    num_weeks integer,
    happiness happiness
);
INSERT INTO holidays(num_weeks,happiness) VALUES (4, 'happy');
INSERT INTO holidays(num_weeks,happiness) VALUES (6, 'very happy');
INSERT INTO holidays(num_weeks,happiness) VALUES (8, 'ecstatic');
INSERT INTO holidays(num_weeks,happiness) VALUES (2, 'sad');
ERROR:  invalid input value for enum happiness: "sad"
SELECT person.name, holidays.num_weeks FROM person, holidays
    WHERE person.current_mood = holidays.happiness;
ERROR:  operator does not exist: mood = happiness
```

If you really need to do something like that, you can either write a custom operator or add explicit casts to your query:

```
SELECT person.name, holidays.num_weeks FROM person, holidays
    WHERE person.current_mood::text = holidays.happiness::text;
name | num_weeks
-----+-----
Moe  |          4
(1 row)
```

8.7.4. Implementation Details

Enum labels are case sensitive, so 'happy' is not the same as 'HAPPY'. White space in the labels is significant too.

Although enum types are primarily intended for static sets of values, there is support for adding new values to an existing enum type, and for renaming values (see [ALTER TYPE](#)). Existing values cannot be removed from an enum type, nor can the sort ordering of such values be changed, short of dropping and re-creating the enum type.

An enum value occupies four bytes on disk. The length of an enum value's textual label is limited by the `NAMEDATALEN` setting compiled into Postgres Pro; in standard builds this means at most 63 bytes.

The translations from internal enum values to textual labels are kept in the system catalog `pg_enum`. Querying this catalog directly can be useful.

8.8. Geometric Types

Geometric data types represent two-dimensional spatial objects. [Table 8.20](#) shows the geometric types available in Postgres Pro.

Table 8.20. Geometric Types

Name	Storage Size	Description	Representation
point	16 bytes	Point on a plane	(x,y)
line	24 bytes	Infinite line	{A,B,C}
lseg	32 bytes	Finite line segment	[(x1,y1),(x2,y2)]
box	32 bytes	Rectangular box	(x1,y1),(x2,y2)
path	16+16n bytes	Closed path (similar to polygon)	((x1,y1),...)

Name	Storage Size	Description	Representation
path	16+16n bytes	Open path	[(x1,y1),...]
polygon	40+16n bytes	Polygon (similar to closed path)	((x1,y1),...)
circle	24 bytes	Circle	<(x,y),r> (center point and radius)

In all these types, the individual coordinates are stored as `double precision (float8)` numbers.

A rich set of functions and operators is available to perform various geometric operations such as scaling, translation, rotation, and determining intersections. They are explained in [Section 9.11](#).

8.8.1. Points

Points are the fundamental two-dimensional building block for geometric types. Values of type `point` are specified using either of the following syntaxes:

```
( x , y )
x , y
```

where `x` and `y` are the respective coordinates, as floating-point numbers.

Points are output using the first syntax.

8.8.2. Lines

Lines are represented by the linear equation $Ax + By + C = 0$, where `A` and `B` are not both zero. Values of type `line` are input and output in the following form:

```
{ A, B, C }
```

Alternatively, any of the following forms can be used for input:

```
[ ( x1 , y1 ) , ( x2 , y2 ) ]
( ( x1 , y1 ) , ( x2 , y2 ) )
( x1 , y1 ) , ( x2 , y2 )
x1 , y1 , x2 , y2
```

where $(x1, y1)$ and $(x2, y2)$ are two different points on the line.

8.8.3. Line Segments

Line segments are represented by pairs of points that are the endpoints of the segment. Values of type `lseg` are specified using any of the following syntaxes:

```
[ ( x1 , y1 ) , ( x2 , y2 ) ]
( ( x1 , y1 ) , ( x2 , y2 ) )
( x1 , y1 ) , ( x2 , y2 )
x1 , y1 , x2 , y2
```

where $(x1, y1)$ and $(x2, y2)$ are the end points of the line segment.

Line segments are output using the first syntax.

8.8.4. Boxes

Boxes are represented by pairs of points that are opposite corners of the box. Values of type `box` are specified using any of the following syntaxes:

```
( ( x1 , y1 ) , ( x2 , y2 ) )
( x1 , y1 ) , ( x2 , y2 )
x1 , y1 , x2 , y2
```

where $(x1, y1)$ and $(x2, y2)$ are any two opposite corners of the box.

Boxes are output using the second syntax.

Any two opposite corners can be supplied on input, but the values will be reordered as needed to store the upper right and lower left corners, in that order.

8.8.5. Paths

Paths are represented by lists of connected points. Paths can be *open*, where the first and last points in the list are considered not connected, or *closed*, where the first and last points are considered connected.

Values of type `path` are specified using any of the following syntaxes:

```
[ ( x1 , y1 ) , ... , ( xn , yn ) ]  
( ( x1 , y1 ) , ... , ( xn , yn ) )  
  ( x1 , y1 ) , ... , ( xn , yn )  
  ( x1 , y1 , ... , xn , yn )  
    x1 , y1 , ... , xn , yn
```

where the points are the end points of the line segments comprising the path. Square brackets (`[]`) indicate an open path, while parentheses (`()`) indicate a closed path. When the outermost parentheses are omitted, as in the third through fifth syntaxes, a closed path is assumed.

Paths are output using the first or second syntax, as appropriate.

8.8.6. Polygons

Polygons are represented by lists of points (the vertices of the polygon). Polygons are very similar to closed paths; the essential semantic difference is that a polygon is considered to include the area within it, while a path is not.

An important implementation difference between polygons and paths is that the stored representation of a polygon includes its smallest bounding box. This speeds up certain search operations, although computing the bounding box adds overhead while constructing new polygons.

Values of type `polygon` are specified using any of the following syntaxes:

```
( ( x1 , y1 ) , ... , ( xn , yn ) )  
  ( x1 , y1 ) , ... , ( xn , yn )  
  ( x1 , y1 , ... , xn , yn )  
    x1 , y1 , ... , xn , yn
```

where the points are the end points of the line segments comprising the boundary of the polygon.

Polygons are output using the first syntax.

8.8.7. Circles

Circles are represented by a center point and radius. Values of type `circle` are specified using any of the following syntaxes:

```
< ( x , y ) , r >  
( ( x , y ) , r )  
  ( x , y ) , r  
    x , y , r
```

where (x, y) is the center point and r is the radius of the circle.

Circles are output using the first syntax.

8.9. Network Address Types

Postgres Pro offers data types to store IPv4, IPv6, and MAC addresses, as shown in [Table 8.21](#). It is better to use these types instead of plain text types to store network addresses, because these types offer input error checking and specialized operators and functions (see [Section 9.12](#)).

Table 8.21. Network Address Types

Name	Storage Size	Description
<code>cidr</code>	7 or 19 bytes	IPv4 and IPv6 networks
<code>inet</code>	7 or 19 bytes	IPv4 and IPv6 hosts and networks
<code>macaddr</code>	6 bytes	MAC addresses
<code>macaddr8</code>	8 bytes	MAC addresses (EUI-64 format)

When sorting `inet` or `cidr` data types, IPv4 addresses will always sort before IPv6 addresses, including IPv4 addresses encapsulated or mapped to IPv6 addresses, such as `::10.2.3.4` or `::ffff:10.4.3.2`.

8.9.1. `inet`

The `inet` type holds an IPv4 or IPv6 host address, and optionally its subnet, all in one field. The subnet is represented by the number of network address bits present in the host address (the “netmask”). If the netmask is 32 and the address is IPv4, then the value does not indicate a subnet, only a single host. In IPv6, the address length is 128 bits, so 128 bits specify a unique host address. Note that if you want to accept only networks, you should use the `cidr` type rather than `inet`.

The input format for this type is `address/y` where `address` is an IPv4 or IPv6 address and `y` is the number of bits in the netmask. If the `/y` portion is omitted, the netmask is taken to be 32 for IPv4 or 128 for IPv6, so the value represents just a single host. On display, the `/y` portion is suppressed if the netmask specifies a single host.

8.9.2. `cidr`

The `cidr` type holds an IPv4 or IPv6 network specification. Input and output formats follow Classless Internet Domain Routing conventions. The format for specifying networks is `address/y` where `address` is the network's lowest address represented as an IPv4 or IPv6 address, and `y` is the number of bits in the netmask. If `y` is omitted, it is calculated using assumptions from the older classful network numbering system, except it will be at least large enough to include all of the octets written in the input. It is an error to specify a network address that has bits set to the right of the specified netmask.

Table 8.22 shows some examples.

Table 8.22. `cidr` Type Input Examples

<code>cidr</code> Input	<code>cidr</code> Output	<code>abbrev(cidr)</code>
192.168.100.128/25	192.168.100.128/25	192.168.100.128/25
192.168/24	192.168.0.0/24	192.168.0/24
192.168/25	192.168.0.0/25	192.168.0.0/25
192.168.1	192.168.1.0/24	192.168.1/24
192.168	192.168.0.0/24	192.168.0/24
128.1	128.1.0.0/16	128.1/16
128	128.0.0.0/16	128.0/16
128.1.2	128.1.2.0/24	128.1.2/24
10.1.2	10.1.2.0/24	10.1.2/24
10.1	10.1.0.0/16	10.1/16
10	10.0.0.0/8	10/8
10.1.2.3/32	10.1.2.3/32	10.1.2.3/32
2001:4f8:3:ba::/64	2001:4f8:3:ba::/64	2001:4f8:3:ba/64
2001:4f8:3:ba:2e0:81f-f:fe22:d1f1/128	2001:4f8:3:ba:2e0:81f-f:fe22:d1f1/128	2001:4f8:3:ba:2e0:81f-f:fe22:d1f1/128

cidr Input	cidr Output	abbrev(cidr)
::ffff:1.2.3.0/120	::ffff:1.2.3.0/120	::ffff:1.2.3/120
::ffff:1.2.3.0/128	::ffff:1.2.3.0/128	::ffff:1.2.3.0/128

8.9.3. inet vs. cidr

The essential difference between `inet` and `cidr` data types is that `inet` accepts values with nonzero bits to the right of the netmask, whereas `cidr` does not. For example, `192.168.0.1/24` is valid for `inet` but not for `cidr`.

Tip

If you do not like the output format for `inet` or `cidr` values, try the functions `host`, `text`, and `abbrev`.

8.9.4. macaddr

The `macaddr` type stores MAC addresses, known for example from Ethernet card hardware addresses (although MAC addresses are used for other purposes as well). Input is accepted in the following formats:

```
'08:00:2b:01:02:03'
'08-00-2b-01-02-03'
'08002b:010203'
'08002b-010203'
'0800.2b01.0203'
'0800-2b01-0203'
'08002b010203'
```

These examples all specify the same address. Upper and lower case is accepted for the digits `a` through `f`. Output is always in the first of the forms shown.

IEEE Standard 802-2001 specifies the second form shown (with hyphens) as the canonical form for MAC addresses, and specifies the first form (with colons) as used with bit-reversed, MSB-first notation, so that `08-00-2b-01-02-03` = `10:00:D4:80:40:C0`. This convention is widely ignored nowadays, and it is relevant only for obsolete network protocols (such as Token Ring). Postgres Pro makes no provisions for bit reversal; all accepted formats use the canonical LSB order.

The remaining five input formats are not part of any standard.

8.9.5. macaddr8

The `macaddr8` type stores MAC addresses in EUI-64 format, known for example from Ethernet card hardware addresses (although MAC addresses are used for other purposes as well). This type can accept both 6 and 8 byte length MAC addresses and stores them in 8 byte length format. MAC addresses given in 6 byte format will be stored in 8 byte length format with the 4th and 5th bytes set to `FF` and `FE`, respectively. Note that IPv6 uses a modified EUI-64 format where the 7th bit should be set to one after the conversion from EUI-48. The function `macaddr8_set7bit` is provided to make this change. Generally speaking, any input which is comprised of pairs of hex digits (on byte boundaries), optionally separated consistently by one of `:`, `-` or `.`, is accepted. The number of hex digits must be either 16 (8 bytes) or 12 (6 bytes). Leading and trailing whitespace is ignored. The following are examples of input formats that are accepted:

```
'08:00:2b:01:02:03:04:05'
'08-00-2b-01-02-03-04-05'
'08002b:0102030405'
'08002b-0102030405'
'0800.2b01.0203.0405'
```

```
'0800-2b01-0203-0405'
'08002b01:02030405'
'08002b0102030405'
```

These examples all specify the same address. Upper and lower case is accepted for the digits a through f. Output is always in the first of the forms shown.

The last six input formats shown above are not part of any standard.

To convert a traditional 48 bit MAC address in EUI-48 format to modified EUI-64 format to be included as the host portion of an IPv6 address, use `macaddr8_set7bit` as shown:

```
SELECT macaddr8_set7bit('08:00:2b:01:02:03');
```

```
      macaddr8_set7bit
-----
0a:00:2b:ff:fe:01:02:03
(1 row)
```

8.10. Bit String Types

Bit strings are strings of 1's and 0's. They can be used to store or visualize bit masks. There are two SQL bit types: `bit(n)` and `bit varying(n)`, where *n* is a positive integer.

`bit` type data must match the length *n* exactly; it is an error to attempt to store shorter or longer bit strings. `bit varying` data is of variable length up to the maximum length *n*; longer strings will be rejected. Writing `bit` without a length is equivalent to `bit(1)`, while `bit varying` without a length specification means unlimited length.

Note

If one explicitly casts a bit-string value to `bit(n)`, it will be truncated or zero-padded on the right to be exactly *n* bits, without raising an error. Similarly, if one explicitly casts a bit-string value to `bit varying(n)`, it will be truncated on the right if it is more than *n* bits.

Refer to [Section 4.1.2.5](#) for information about the syntax of bit string constants. Bit-logical operators and string manipulation functions are available; see [Section 9.6](#).

Example 8.3. Using the Bit String Types

```
CREATE TABLE test (a BIT(3), b BIT VARYING(5));
INSERT INTO test VALUES (B'101', B'00');
INSERT INTO test VALUES (B'10', B'101');
```

```
ERROR: bit string length 2 does not match type bit(3)
```

```
INSERT INTO test VALUES (B'10'::bit(3), B'101');
SELECT * FROM test;
```

```
  a | b
-----+-----
101 | 00
100 | 101
```

A bit string value requires 1 byte for each group of 8 bits, plus 5 or 8 bytes overhead depending on the length of the string (but long values may be compressed or moved out-of-line, as explained in [Section 8.3](#) for character strings).

8.11. Text Search Types

Postgres Pro provides two data types that are designed to support full text search, which is the activity of searching through a collection of natural-language *documents* to locate those that best match a *query*. The `tsvector` type represents a document in a form optimized for text search; the `tsquery` type similarly represents a text query. [Chapter 12](#) provides a detailed explanation of this facility, and [Section 9.13](#) summarizes the related functions and operators.

8.11.1. `tsvector`

A `tsvector` value is a sorted list of distinct *lexemes*, which are words that have been *normalized* to merge different variants of the same word (see [Chapter 12](#) for details). Sorting and duplicate-elimination are done automatically during input, as shown in this example:

```
SELECT 'a fat cat sat on a mat and ate a fat rat'::tsvector;
           tsvector
-----
'a' 'and' 'ate' 'cat' 'fat' 'mat' 'on' 'rat' 'sat'
```

To represent lexemes containing whitespace or punctuation, surround them with quotes:

```
SELECT $$the lexeme '    ' contains spaces$$::tsvector;
           tsvector
-----
'    ' 'contains' 'lexeme' 'spaces' 'the'
```

(We use dollar-quoted string literals in this example and the next one to avoid the confusion of having to double quote marks within the literals.) Embedded quotes and backslashes must be doubled:

```
SELECT $$the lexeme 'Joe''s' contains a quote$$::tsvector;
           tsvector
-----
'Joe''s' 'a' 'contains' 'lexeme' 'quote' 'the'
```

Optionally, integer *positions* can be attached to lexemes:

```
SELECT 'a:1 fat:2 cat:3 sat:4 on:5 a:6 mat:7 and:8 ate:9 a:10 fat:11 rat:12'::tsvector;
           tsvector
-----
'a':1,6,10 'and':8 'ate':9 'cat':3 'fat':2,11 'mat':7 'on':5 'rat':12 'sat':4
```

A position normally indicates the source word's location in the document. Positional information can be used for *proximity ranking*. Position values can range from 1 to 16383; larger numbers are silently set to 16383. Duplicate positions for the same lexeme are discarded.

Lexemes that have positions can further be labeled with a *weight*, which can be A, B, C, or D. D is the default and hence is not shown on output:

```
SELECT 'a:1A fat:2B,4C cat:5D'::tsvector;
           tsvector
-----
'a':1A 'cat':5 'fat':2B,4C
```

Weights are typically used to reflect document structure, for example by marking title words differently from body words. Text search ranking functions can assign different priorities to the different weight markers.

It is important to understand that the `tsvector` type itself does not perform any word normalization; it assumes the words it is given are normalized appropriately for the application. For example,

```
SELECT 'The Fat Rats'::tsvector;
           tsvector
-----
'Fat' 'Rats' 'The'
```

For most English-text-searching applications the above words would be considered non-normalized, but `tsvector` doesn't care. Raw document text should usually be passed through `to_tsvector` to normalize the words appropriately for searching:

```
SELECT to_tsvector('english', 'The Fat Rats');
       to_tsvector
-----
'fat':2 'rat':3
```

Again, see [Chapter 12](#) for more detail.

8.11.2. tsquery

A `tsquery` value stores lexemes that are to be searched for, and can combine them using the Boolean operators `&` (AND), `|` (OR), and `!` (NOT), as well as the phrase search operator `<->` (FOLLOWED BY). There is also a variant `<N>` of the FOLLOWED BY operator, where *N* is an integer constant that specifies the distance between the two lexemes being searched for. `<->` is equivalent to `<1>`.

Parentheses can be used to enforce grouping of these operators. In the absence of parentheses, `!` (NOT) binds most tightly, `<->` (FOLLOWED BY) next most tightly, then `&` (AND), with `|` (OR) binding the least tightly.

Here are some examples:

```
SELECT 'fat & rat'::tsquery;
       tsquery
-----
'fat' & 'rat'

SELECT 'fat & (rat | cat)'::tsquery;
       tsquery
-----
'fat' & ( 'rat' | 'cat' )

SELECT 'fat & rat & ! cat'::tsquery;
       tsquery
-----
'fat' & 'rat' & !'cat'
```

Optionally, lexemes in a `tsquery` can be labeled with one or more weight letters, which restricts them to match only `tsvector` lexemes with one of those weights:

```
SELECT 'fat:ab & cat'::tsquery;
       tsquery
-----
'fat':AB & 'cat'
```

Also, lexemes in a `tsquery` can be labeled with `*` to specify prefix matching:

```
SELECT 'super:*'::tsquery;
       tsquery
-----
'super':*
```

This query will match any word in a `tsvector` that begins with “super”.

Quoting rules for lexemes are the same as described previously for lexemes in `tsvector`; and, as with `tsvector`, any required normalization of words must be done before converting to the `tsquery` type. The `to_tsquery` function is convenient for performing such normalization:

```
SELECT to_tsquery('Fat:ab & Cats');
       to_tsquery
```

```
-----  
'fat':AB & 'cat'
```

Note that `to_tsquery` will process prefixes in the same way as other words, which means this comparison returns true:

```
SELECT to_tsvector( 'postgraduate' ) @@ to_tsquery( 'postgres:*' );  
?column?
```

```
-----  
t
```

because `postgres` gets stemmed to `postgr`:

```
SELECT to_tsvector( 'postgraduate' ), to_tsquery( 'postgres:*' );  
to_tsvector | to_tsquery  
-----+-----  
'postgradu':1 | 'postgr':*
```

which will match the stemmed form of `postgraduate`.

8.12. UUID Type

The data type `uuid` stores Universally Unique Identifiers (UUID) as defined by [RFC 4122](#), ISO/IEC 9834-8:2005, and related standards. (Some systems refer to this data type as a globally unique identifier, or GUID, instead.) This identifier is a 128-bit quantity that is generated by an algorithm chosen to make it very unlikely that the same identifier will be generated by anyone else in the known universe using the same algorithm. Therefore, for distributed systems, these identifiers provide a better uniqueness guarantee than sequence generators, which are only unique within a single database.

A UUID is written as a sequence of lower-case hexadecimal digits, in several groups separated by hyphens, specifically a group of 8 digits followed by three groups of 4 digits followed by a group of 12 digits, for a total of 32 digits representing the 128 bits. An example of a UUID in this standard form is:

```
a0eebc99-9c0b-4ef8-bb6d-6bb9bd380a11
```

Postgres Pro also accepts the following alternative forms for input: use of upper-case digits, the standard format surrounded by braces, omitting some or all hyphens, adding a hyphen after any group of four digits. Examples are:

```
A0EEBC99-9C0B-4EF8-BB6D-6BB9BD380A11  
{a0eebc99-9c0b-4ef8-bb6d-6bb9bd380a11}  
a0eebc999c0b4ef8bb6d6bb9bd380a11  
a0ee-bc99-9c0b-4ef8-bb6d-6bb9-bd38-0a11  
{a0eebc99-9c0b4ef8-bb6d6bb9-bd380a11}
```

Output is always in the standard form.

See [Section 9.14](#) for how to generate a UUID in Postgres Pro.

8.13. XML Type

The `xml` data type can be used to store XML data. Its advantage over storing XML data in a `text` field is that it checks the input values for well-formedness, and there are support functions to perform type-safe operations on it; see [Section 9.15](#). Use of this data type requires the installation to have been built with `configure --with-libxml`.

The `xml` type can store well-formed “documents”, as defined by the XML standard, as well as “content” fragments, which are defined by reference to the more permissive “*document node*” of the XQuery and XPath data model. Roughly, this means that content fragments can have more than one top-level element or character node. The expression `xmlvalue IS DOCUMENT` can be used to evaluate whether a particular `xml` value is a full document or only a content fragment.

Limits and compatibility notes for the `xml` data type can be found in [Section D.3](#).

8.13.1. Creating XML Values

To produce a value of type `xml` from character data, use the function `xmlparse`:

```
XMLPARSE ( { DOCUMENT | CONTENT } value)
```

Examples:

```
XMLPARSE (DOCUMENT '<?xml version="1.0"?><book><title>Manual</title><chapter>...</chapter></book>')
XMLPARSE (CONTENT 'abc<foo>bar</foo><bar>foo</bar>')
```

While this is the only way to convert character strings into XML values according to the SQL standard, the Postgres Pro-specific syntaxes:

```
xml '<foo>bar</foo>'
'<foo>bar</foo>'::xml
```

can also be used.

The `xml` type does not validate input values against a document type declaration (DTD), even when the input value specifies a DTD. There is also currently no built-in support for validating against other XML schema languages such as XML Schema.

The inverse operation, producing a character string value from `xml`, uses the function `xmlserialize`:

```
XMLSERIALIZE ( { DOCUMENT | CONTENT } value AS type [ [ NO ] INDENT ] )
```

`type` can be `character`, `character varying`, or `text` (or an alias for one of those). Again, according to the SQL standard, this is the only way to convert between type `xml` and character types, but Postgres Pro also allows you to simply cast the value.

The `INDENT` option causes the result to be pretty-printed, while `NO INDENT` (which is the default) just emits the original input string. Casting to a character type likewise produces the original string.

When a character string value is cast to or from type `xml` without going through `XMLPARSE` or `XMLSERIALIZE`, respectively, the choice of `DOCUMENT` versus `CONTENT` is determined by the “XML option” session configuration parameter, which can be set using the standard command:

```
SET XML OPTION { DOCUMENT | CONTENT };
```

or the more Postgres Pro-like syntax

```
SET xmloption TO { DOCUMENT | CONTENT };
```

The default is `CONTENT`, so all forms of XML data are allowed.

8.13.2. Encoding Handling

Care must be taken when dealing with multiple character encodings on the client, server, and in the XML data passed through them. When using the text mode to pass queries to the server and query results to the client (which is the normal mode), Postgres Pro converts all character data passed between the client and the server and vice versa to the character encoding of the respective end; see [Section 23.3](#). This includes string representations of XML values, such as in the above examples. This would ordinarily mean that encoding declarations contained in XML data can become invalid as the character data is converted to other encodings while traveling between client and server, because the embedded encoding declaration is not changed. To cope with this behavior, encoding declarations contained in character strings presented for input to the `xml` type are *ignored*, and content is assumed to be in the current server encoding. Consequently, for correct processing, character strings of XML data must be sent from the client in the current client encoding. It is the responsibility of the client to either convert documents to the current client encoding before sending them to the server, or to adjust the client encoding appropriately. On output, values of type `xml` will not have an encoding declaration, and clients should assume all data is in the current client encoding.

When using binary mode to pass query parameters to the server and query results back to the client, no encoding conversion is performed, so the situation is different. In this case, an encoding declaration in

the XML data will be observed, and if it is absent, the data will be assumed to be in UTF-8 (as required by the XML standard; note that Postgres Pro does not support UTF-16). On output, data will have an encoding declaration specifying the client encoding, unless the client encoding is UTF-8, in which case it will be omitted.

Needless to say, processing XML data with Postgres Pro will be less error-prone and more efficient if the XML data encoding, client encoding, and server encoding are the same. Since XML data is internally processed in UTF-8, computations will be most efficient if the server encoding is also UTF-8.

Caution

Some XML-related functions may not work at all on non-ASCII data when the server encoding is not UTF-8. This is known to be an issue for `xmltable()` and `xpath()` in particular.

8.13.3. Accessing XML Values

The `xml` data type is unusual in that it does not provide any comparison operators. This is because there is no well-defined and universally useful comparison algorithm for XML data. One consequence of this is that you cannot retrieve rows by comparing an `xml` column against a search value. XML values should therefore typically be accompanied by a separate key field such as an ID. An alternative solution for comparing XML values is to convert them to character strings first, but note that character string comparison has little to do with a useful XML comparison method.

Since there are no comparison operators for the `xml` data type, it is not possible to create an index directly on a column of this type. If speedy searches in XML data are desired, possible workarounds include casting the expression to a character string type and indexing that, or indexing an XPath expression. Of course, the actual query would have to be adjusted to search by the indexed expression.

The text-search functionality in Postgres Pro can also be used to speed up full-document searches of XML data. The necessary preprocessing support is, however, not yet available in the Postgres Pro distribution.

8.14. JSON Types

JSON data types are for storing JSON (JavaScript Object Notation) data, as specified in [RFC 7159](#). Such data can also be stored as `text`, but the JSON data types have the advantage of enforcing that each stored value is valid according to the JSON rules. There are also assorted JSON-specific functions and operators available for data stored in these data types; see [Section 9.16](#).

Postgres Pro offers two types for storing JSON data: `json` and `jsonb`. To implement efficient query mechanisms for these data types, Postgres Pro also provides the `jsonpath` data type described in [Section 8.14.7](#).

The `json` and `jsonb` data types accept *almost* identical sets of values as input. The major practical difference is one of efficiency. The `json` data type stores an exact copy of the input text, which processing functions must reparse on each execution; while `jsonb` data is stored in a decomposed binary format that makes it slightly slower to input due to added conversion overhead, but significantly faster to process, since no reparsing is needed. `jsonb` also supports indexing, which can be a significant advantage.

Because the `json` type stores an exact copy of the input text, it will preserve semantically-insignificant white space between tokens, as well as the order of keys within JSON objects. Also, if a JSON object within the value contains the same key more than once, all the key/value pairs are kept. (The processing functions consider the last value as the operative one.) By contrast, `jsonb` does not preserve white space, does not preserve the order of object keys, and does not keep duplicate object keys. If duplicate keys are specified in the input, only the last value is kept.

In general, most applications should prefer to store JSON data as `jsonb`, unless there are quite specialized needs, such as legacy assumptions about ordering of object keys.

RFC 7159 specifies that JSON strings should be encoded in UTF8. It is therefore not possible for the JSON types to conform rigidly to the JSON specification unless the database encoding is UTF8. Attempts to directly include characters that cannot be represented in the database encoding will fail; conversely, characters that can be represented in the database encoding but not in UTF8 will be allowed.

RFC 7159 permits JSON strings to contain Unicode escape sequences denoted by `\uXXXX`. In the input function for the `json` type, Unicode escapes are allowed regardless of the database encoding, and are checked only for syntactic correctness (that is, that four hex digits follow `\u`). However, the input function for `jsonb` is stricter: it disallows Unicode escapes for characters that cannot be represented in the database encoding. The `jsonb` type also rejects `\u0000` (because that cannot be represented in Postgres Pro's `text` type), and it insists that any use of Unicode surrogate pairs to designate characters outside the Unicode Basic Multilingual Plane be correct. Valid Unicode escapes are converted to the equivalent single character for storage; this includes folding surrogate pairs into a single character.

Note

Many of the JSON processing functions described in [Section 9.16](#) will convert Unicode escapes to regular characters, and will therefore throw the same types of errors just described even if their input is of type `json` not `jsonb`. The fact that the `json` input function does not make these checks may be considered a historical artifact, although it does allow for simple storage (without processing) of JSON Unicode escapes in a database encoding that does not support the represented characters.

When converting textual JSON input into `jsonb`, the primitive types described by RFC 7159 are effectively mapped onto native Postgres Pro types, as shown in [Table 8.23](#). Therefore, there are some minor additional constraints on what constitutes valid `jsonb` data that do not apply to the `json` type, nor to JSON in the abstract, corresponding to limits on what can be represented by the underlying data type. Notably, `jsonb` will reject numbers that are outside the range of the Postgres Pro `numeric` data type, while `json` will not. Such implementation-defined restrictions are permitted by RFC 7159. However, in practice such problems are far more likely to occur in other implementations, as it is common to represent JSON's `number` primitive type as IEEE 754 double precision floating point (which RFC 7159 explicitly anticipates and allows for). When using JSON as an interchange format with such systems, the danger of losing numeric precision compared to data originally stored by Postgres Pro should be considered.

Conversely, as noted in the table there are some minor restrictions on the input format of JSON primitive types that do not apply to the corresponding Postgres Pro types.

Table 8.23. JSON Primitive Types and Corresponding Postgres Pro Types

JSON primitive type	Postgres Pro type	Notes
string	text	<code>\u0000</code> is disallowed, as are Unicode escapes representing characters not available in the database encoding
number	numeric	NaN and infinity values are disallowed
boolean	boolean	Only lowercase <code>true</code> and <code>false</code> spellings are accepted
null	(none)	SQL NULL is a different concept

8.14.1. JSON Input and Output Syntax

The input/output syntax for the JSON data types is as specified in RFC 7159.

The following are all valid `json` (or `jsonb`) expressions:

```
-- Simple scalar/primitive value
-- Primitive values can be numbers, quoted strings, true, false, or null
SELECT '5'::json;
```

```
-- Array of zero or more elements (elements need not be of same type)
SELECT '[1, 2, "foo", null]':::json;

-- Object containing pairs of keys and values
-- Note that object keys must always be quoted strings
SELECT '{"bar": "baz", "balance": 7.77, "active": false}':::json;

-- Arrays and objects can be nested arbitrarily
SELECT '{"foo": [true, "bar"], "tags": {"a": 1, "b": null}}':::json;
```

As previously stated, when a JSON value is input and then printed without any additional processing, `json` outputs the same text that was input, while `jsonb` does not preserve semantically-insignificant details such as whitespace. For example, note the differences here:

```
SELECT '{"bar": "baz", "balance": 7.77, "active":false}':::json;
               json
-----
{"bar": "baz", "balance": 7.77, "active":false}
(1 row)

SELECT '{"bar": "baz", "balance": 7.77, "active":false}':::jsonb;
               jsonb
-----
{"bar": "baz", "active": false, "balance": 7.77}
(1 row)
```

One semantically-insignificant detail worth noting is that in `jsonb`, numbers will be printed according to the behavior of the underlying `numeric` type. In practice this means that numbers entered with `E` notation will be printed without it, for example:

```
SELECT '{"reading": 1.230e-5}':::json, '{"reading": 1.230e-5}':::jsonb;
               json               |               jsonb
-----+-----
{"reading": 1.230e-5} | {"reading": 0.00001230}
(1 row)
```

However, `jsonb` will preserve trailing fractional zeroes, as seen in this example, even though those are semantically insignificant for purposes such as equality checks.

For the list of built-in functions and operators available for constructing and processing JSON values, see [Section 9.16](#).

8.14.2. Designing JSON Documents

Representing data as JSON can be considerably more flexible than the traditional relational data model, which is compelling in environments where requirements are fluid. It is quite possible for both approaches to co-exist and complement each other within the same application. However, even for applications where maximal flexibility is desired, it is still recommended that JSON documents have a somewhat fixed structure. The structure is typically unenforced (though enforcing some business rules declaratively is possible), but having a predictable structure makes it easier to write queries that usefully summarize a set of “documents” (datums) in a table.

JSON data is subject to the same concurrency-control considerations as any other data type when stored in a table. Although storing large documents is practicable, keep in mind that any update acquires a row-level lock on the whole row. Consider limiting JSON documents to a manageable size in order to decrease lock contention among updating transactions. Ideally, JSON documents should each represent an atomic datum that business rules dictate cannot reasonably be further subdivided into smaller datums that could be modified independently.

8.14.3. `jsonb` Containment and Existence

Testing *containment* is an important capability of `jsonb`. There is no parallel set of facilities for the `json` type. Containment tests whether one `jsonb` document has contained within it another one. These examples return true except as noted:

```
-- Simple scalar/primitive values contain only the identical value:
SELECT '"foo"'::jsonb @> '"foo"'::jsonb;

-- The array on the right side is contained within the one on the left:
SELECT '[1, 2, 3]'::jsonb @> '[1, 3]'::jsonb;

-- Order of array elements is not significant, so this is also true:
SELECT '[1, 2, 3]'::jsonb @> '[3, 1]'::jsonb;

-- Duplicate array elements don't matter either:
SELECT '[1, 2, 3]'::jsonb @> '[1, 2, 2]'::jsonb;

-- The object with a single pair on the right side is contained
-- within the object on the left side:
SELECT '{"product": "PostgreSQL", "version": 9.4, "jsonb": true}'::jsonb @>
'{"version": 9.4}'::jsonb;

-- The array on the right side is not considered contained within the
-- array on the left, even though a similar array is nested within it:
SELECT '[1, 2, [1, 3]]'::jsonb @> '[1, 3]'::jsonb; -- yields false

-- But with a layer of nesting, it is contained:
SELECT '[1, 2, [1, 3]]'::jsonb @> '[[1, 3]]'::jsonb;

-- Similarly, containment is not reported here:
SELECT '{"foo": {"bar": "baz"}}'::jsonb @> '{"bar": "baz"}'::jsonb; -- yields false

-- A top-level key and an empty object is contained:
SELECT '{"foo": {"bar": "baz"}}'::jsonb @> '{"foo": {}}'::jsonb;
```

The general principle is that the contained object must match the containing object as to structure and data contents, possibly after discarding some non-matching array elements or object key/value pairs from the containing object. But remember that the order of array elements is not significant when doing a containment match, and duplicate array elements are effectively considered only once.

As a special exception to the general principle that the structures must match, an array may contain a primitive value:

```
-- This array contains the primitive string value:
SELECT '["foo", "bar"]'::jsonb @> '"bar"'::jsonb;

-- This exception is not reciprocal -- non-containment is reported here:
SELECT '"bar"'::jsonb @> '["bar"]'::jsonb; -- yields false
```

`jsonb` also has an *existence* operator, which is a variation on the theme of containment: it tests whether a string (given as a text value) appears as an object key or array element at the top level of the `jsonb` value. These examples return true except as noted:

```
-- String exists as array element:
SELECT '["foo", "bar", "baz"]'::jsonb ? 'bar';

-- String exists as object key:
SELECT '{"foo": "bar"}'::jsonb ? 'foo';

-- Object values are not considered:
```



```
SELECT '{"foo": "bar"}'::jsonb ? 'bar'; -- yields false

-- As with containment, existence must match at the top level:
SELECT '{"foo": {"bar": "baz"}}'::jsonb ? 'bar'; -- yields false

-- A string is considered to exist if it matches a primitive JSON string:
SELECT '"foo"'::jsonb ? 'foo';
```

JSON objects are better suited than arrays for testing containment or existence when there are many keys or elements involved, because unlike arrays they are internally optimized for searching, and do not need to be searched linearly.

Tip

Because JSON containment is nested, an appropriate query can skip explicit selection of sub-objects. As an example, suppose that we have a `doc` column containing objects at the top level, with most objects containing `tags` fields that contain arrays of sub-objects. This query finds entries in which sub-objects containing both `"term": "paris"` and `"term": "food"` appear, while ignoring any such keys outside the `tags` array:

```
SELECT doc->'site_name' FROM websites
WHERE doc @> '{"tags":[{"term":"paris"}, {"term":"food"}]}';
```

One could accomplish the same thing with, say,

```
SELECT doc->'site_name' FROM websites
WHERE doc->'tags' @> '[{"term":"paris"}, {"term":"food"}]';
```

but that approach is less flexible, and often less efficient as well.

On the other hand, the JSON existence operator is not nested: it will only look for the specified key or array element at top level of the JSON value.

The various containment and existence operators, along with all other JSON operators and functions are documented in [Section 9.16](#).

8.14.4. jsonb Indexing

GIN indexes can be used to efficiently search for keys or key/value pairs occurring within a large number of `jsonb` documents (datums). Two GIN “operator classes” are provided, offering different performance and flexibility trade-offs.

The default GIN operator class for `jsonb` supports queries with the key-exists operators `?`, `?|` and `?&`, the containment operator `@>`, and the `jsonpath` match operators `@?` and `@@`. (For details of the semantics that these operators implement, see [Table 9.46](#).) An example of creating an index with this operator class is:

```
CREATE INDEX idxgin ON api USING GIN (jdoc);
```

The non-default GIN operator class `jsonb_path_ops` does not support the key-exists operators, but it does support `@>`, `@?` and `@@`. An example of creating an index with this operator class is:

```
CREATE INDEX idxginp ON api USING GIN (jdoc jsonb_path_ops);
```

Consider the example of a table that stores JSON documents retrieved from a third-party web service, with a documented schema definition. A typical document is:

```
{
  "guid": "9c36adc1-7fb5-4d5b-83b4-90356a46061a",
  "name": "Angela Barton",
  "is_active": true,
  "company": "MagnaFone",
  "address": "178 Howard Place, Gulf, Washington, 702",
```

```
"registered": "2009-11-07T08:53:22 +08:00",
"latitude": 19.793713,
"longitude": 86.513373,
"tags": [
    "enim",
    "aliquip",
    "qui"
]
}
```

We store these documents in a table named `api`, in a `jsonb` column named `jdoc`. If a GIN index is created on this column, queries like the following can make use of the index:

```
-- Find documents in which the key "company" has value "MagnaFone"
SELECT jdoc->'guid', jdoc->'name' FROM api WHERE jdoc @> '{"company": "MagnaFone"}';
```

However, the index could not be used for queries like the following, because though the operator `?` is indexable, it is not applied directly to the indexed column `jdoc`:

```
-- Find documents in which the key "tags" contains key or array element "qui"
SELECT jdoc->'guid', jdoc->'name' FROM api WHERE jdoc -> 'tags' ? 'qui';
```

Still, with appropriate use of expression indexes, the above query can use an index. If querying for particular items within the `"tags"` key is common, defining an index like this may be worthwhile:

```
CREATE INDEX idxgintags ON api USING GIN ((jdoc -> 'tags'));
```

Now, the `WHERE` clause `jdoc -> 'tags' ? 'qui'` will be recognized as an application of the indexable operator `?` to the indexed expression `jdoc -> 'tags'`. (More information on expression indexes can be found in [Section 11.7](#).)

Another approach to querying is to exploit containment, for example:

```
-- Find documents in which the key "tags" contains array element "qui"
SELECT jdoc->'guid', jdoc->'name' FROM api WHERE jdoc @> '{"tags": ["qui"]}';
```

A simple GIN index on the `jdoc` column can support this query. But note that such an index will store copies of every key and value in the `jdoc` column, whereas the expression index of the previous example stores only data found under the `tags` key. While the simple-index approach is far more flexible (since it supports queries about any key), targeted expression indexes are likely to be smaller and faster to search than a simple index.

GIN indexes also support the `@?` and `@@` operators, which perform `jsonpath` matching. Examples are

```
SELECT jdoc->'guid', jdoc->'name' FROM api WHERE jdoc @? '$.tags[*] ? (@ == "qui")';
SELECT jdoc->'guid', jdoc->'name' FROM api WHERE jdoc @@ '$.tags[*] == "qui";
```

For these operators, a GIN index extracts clauses of the form `accessors_chain == constant` out of the `jsonpath` pattern, and does the index search based on the keys and values mentioned in these clauses. The accessors chain may include `.key`, `[*]`, and `[index]` accessors. The `jsonb_ops` operator class also supports `.*` and `**` accessors, but the `jsonb_path_ops` operator class does not.

Although the `jsonb_path_ops` operator class supports only queries with the `@>`, `@?` and `@@` operators, it has notable performance advantages over the default operator class `jsonb_ops`. A `jsonb_path_ops` index is usually much smaller than a `jsonb_ops` index over the same data, and the specificity of searches is better, particularly when queries contain keys that appear frequently in the data. Therefore search operations typically perform better than with the default operator class.

The technical difference between a `jsonb_ops` and a `jsonb_path_ops` GIN index is that the former creates independent index items for each key and value in the data, while the latter creates index items only for each value in the data.¹ Basically, each `jsonb_path_ops` index item is a hash of the value and the key(s) leading to it; for example to index `{"foo": {"bar": "baz"}}`, a single index item would

¹ For this purpose, the term “value” includes array elements, though JSON terminology sometimes considers array elements distinct from values within objects.

be created incorporating all three of `foo`, `bar`, and `baz` into the hash value. Thus a containment query looking for this structure would result in an extremely specific index search; but there is no way at all to find out whether `foo` appears as a key. On the other hand, a `jsonb_ops` index would create three index items representing `foo`, `bar`, and `baz` separately; then to do the containment query, it would look for rows containing all three of these items. While GIN indexes can perform such an AND search fairly efficiently, it will still be less specific and slower than the equivalent `jsonb_path_ops` search, especially if there are a very large number of rows containing any single one of the three index items.

A disadvantage of the `jsonb_path_ops` approach is that it produces no index entries for JSON structures not containing any values, such as `{"a": {}}`. If a search for documents containing such a structure is requested, it will require a full-index scan, which is quite slow. `jsonb_path_ops` is therefore ill-suited for applications that often perform such searches.

`jsonb` also supports `btree` and `hash` indexes. These are usually useful only if it's important to check equality of complete JSON documents. The `btree` ordering for `jsonb` datums is seldom of great interest, but for completeness it is:

```
Object > Array > Boolean > Number > String > null
```

```
Object with n pairs > object with n - 1 pairs
```

```
Array with n elements > array with n - 1 elements
```

with the exception that (for historical reasons) an empty top level array sorts less than `null`. Objects with equal numbers of pairs are compared in the order:

```
key-1, value-1, key-2 ...
```

Note that object keys are compared in their storage order; in particular, since shorter keys are stored before longer keys, this can lead to results that might be unintuitive, such as:

```
{ "aa": 1, "c": 1 } > { "b": 1, "d": 1 }
```

Similarly, arrays with equal numbers of elements are compared in the order:

```
element-1, element-2 ...
```

Primitive JSON values are compared using the same comparison rules as for the underlying Postgres Pro data type. Strings are compared using the default database collation.

8.14.5. `jsonb` Subscripting

The `jsonb` data type supports array-style subscripting expressions to extract and modify elements. Nested values can be indicated by chaining subscripting expressions, following the same rules as the `path` argument in the `jsonb_set` function. If a `jsonb` value is an array, numeric subscripts start at zero, and negative integers count backwards from the last element of the array. Slice expressions are not supported. The result of a subscripting expression is always of the `jsonb` data type.

`UPDATE` statements may use subscripting in the `SET` clause to modify `jsonb` values. Subscript paths must be traversable for all affected values insofar as they exist. For instance, the path `val['a']['b']['c']` can be traversed all the way to `c` if every `val`, `val['a']`, and `val['a']['b']` is an object. If any `val['a']` or `val['a']['b']` is not defined, it will be created as an empty object and filled as necessary. However, if any `val` itself or one of the intermediary values is defined as a non-object such as a string, number, or `jsonb null`, traversal cannot proceed so an error is raised and the transaction aborted.

An example of subscripting syntax:

```
-- Extract object value by key
SELECT ('{"a": 1}'::jsonb)['a'];

-- Extract nested object value by key path
SELECT ('{"a": {"b": {"c": 1}}}'::jsonb)['a']['b']['c'];
```

```
-- Extract array element by index
SELECT ('[1, "2", null] '::jsonb)[1];

-- Update object value by key. Note the quotes around '1': the assigned
-- value must be of the jsonb type as well
UPDATE table_name SET jsonb_field['key'] = '1';

-- This will raise an error if any record's jsonb_field['a']['b'] is something
-- other than an object. For example, the value {"a": 1} has a numeric value
-- of the key 'a'.
UPDATE table_name SET jsonb_field['a']['b']['c'] = '1';

-- Filter records using a WHERE clause with subscripting. Since the result of
-- subscripting is jsonb, the value we compare it against must also be jsonb.
-- The double quotes make "value" also a valid jsonb string.
SELECT * FROM table_name WHERE jsonb_field['key'] = '"value"';
```

jsonb assignment via subscripting handles a few edge cases differently from jsonb_set. When a source jsonb value is NULL, assignment via subscripting will proceed as if it was an empty JSON value of the type (object or array) implied by the subscript key:

```
-- Where jsonb_field was NULL, it is now {"a": 1}
UPDATE table_name SET jsonb_field['a'] = '1';

-- Where jsonb_field was NULL, it is now [1]
UPDATE table_name SET jsonb_field[0] = '1';
```

If an index is specified for an array containing too few elements, NULL elements will be appended until the index is reachable and the value can be set.

```
-- Where jsonb_field was [], it is now [null, null, 2];
-- where jsonb_field was [0], it is now [0, null, 2]
UPDATE table_name SET jsonb_field[2] = '2';
```

A jsonb value will accept assignments to nonexistent subscript paths as long as the last existing element to be traversed is an object or array, as implied by the corresponding subscript (the element indicated by the last subscript in the path is not traversed and may be anything). Nested array and object structures will be created, and in the former case null-padded, as specified by the subscript path until the assigned value can be placed.

```
-- Where jsonb_field was {}, it is now {"a": [{"b": 1}]}
UPDATE table_name SET jsonb_field['a'][0]['b'] = '1';

-- Where jsonb_field was [], it is now [null, {"a": 1}]
UPDATE table_name SET jsonb_field[1]['a'] = '1';
```

8.14.6. Transforms

Additional extensions are available that implement transforms for the jsonb type for different procedural languages.

The extensions for PL/Perl are called jsonb_plperl and jsonb_plperl_u. If you use them, jsonb values are mapped to Perl arrays, hashes, and scalars, as appropriate.

The extension for PL/Python is called jsonb_plpython3u. If you use it, jsonb values are mapped to Python dictionaries, lists, and scalars, as appropriate.

Of these extensions, jsonb_plperl is considered “trusted”, that is, it can be installed by non-superusers who have CREATE privilege on the current database. The rest require superuser privilege to install.

8.14.7. jsonpath Type

The `jsonpath` type implements support for the SQL/JSON path language in Postgres Pro to efficiently query JSON data. It provides a binary representation of the parsed SQL/JSON path expression that specifies the items to be retrieved by the path engine from the JSON data for further processing with the SQL/JSON query functions.

The semantics of SQL/JSON path predicates and operators generally follow SQL. At the same time, to provide a natural way of working with JSON data, SQL/JSON path syntax uses some JavaScript conventions:

- Dot (`.`) is used for member access.
- Square brackets (`[]`) are used for array access.
- SQL/JSON arrays are 0-relative, unlike regular SQL arrays that start from 1.

Numeric literals in SQL/JSON path expressions follow JavaScript rules, which are different from both SQL and JSON in some minor details. For example, SQL/JSON path allows `.1` and `1.`, which are invalid in JSON. Non-decimal integer literals and underscore separators are supported, for example, `1_000_000`, `0x1EEE_FFFF`, `0o273`, `0b100101`. In SQL/JSON path (and in JavaScript, but not in SQL proper), there must not be an underscore separator directly after the radix prefix.

An SQL/JSON path expression is typically written in an SQL query as an SQL character string literal, so it must be enclosed in single quotes, and any single quotes desired within the value must be doubled (see [Section 4.1.2.1](#)). Some forms of path expressions require string literals within them. These embedded string literals follow JavaScript/ECMAScript conventions: they must be surrounded by double quotes, and backslash escapes may be used within them to represent otherwise-hard-to-type characters. In particular, the way to write a double quote within an embedded string literal is `\`, and to write a backslash itself, you must write `\\`. Other special backslash sequences include those recognized in JavaScript strings: `\b`, `\f`, `\n`, `\r`, `\t`, `\v` for various ASCII control characters, `\xNN` for a character code written with only two hex digits, `\uNNNN` for a Unicode character identified by its 4-hex-digit code point, and `\u{N...}` for a Unicode character code point written with 1 to 6 hex digits.

A path expression consists of a sequence of path elements, which can be any of the following:

- Path literals of JSON primitive types: Unicode text, numeric, true, false, or null.
- Path variables listed in [Table 8.24](#).
- Accessor operators listed in [Table 8.25](#).
- `jsonpath` operators and methods listed in [Section 9.16.2.3](#).
- Parentheses, which can be used to provide filter expressions or define the order of path evaluation.

For details on using `jsonpath` expressions with SQL/JSON query functions, see [Section 9.16.2](#).

Table 8.24. `jsonpath` Variables

Variable	Description
<code>\$</code>	A variable representing the JSON value being queried (the <i>context item</i>).
<code>\$varname</code>	A named variable. Its value can be set by the parameter <code>vars</code> of several JSON processing functions; see Table 9.49 for details.
<code>@</code>	A variable representing the result of path evaluation in filter expressions.

Table 8.25. `jsonpath` Accessors

Accessor Operator	Description
<code>.key</code>	Member accessor that returns an object member with the specified key. If the key name matches some named variable starting with <code>\$</code>
<code>."\$varname"</code>	

Accessor Operator	Description
	or does not meet the JavaScript rules for an identifier, it must be enclosed in double quotes to make it a string literal.
<code>.*</code>	Wildcard member accessor that returns the values of all members located at the top level of the current object.
<code>.**</code>	Recursive wildcard member accessor that processes all levels of the JSON hierarchy of the current object and returns all the member values, regardless of their nesting level. This is a Postgres Pro extension of the SQL/JSON standard.
<code>.**{level}</code> <code>.**{start_level to end_level}</code>	Like <code>.**</code> , but selects only the specified levels of the JSON hierarchy. Nesting levels are specified as integers. Level zero corresponds to the current object. To access the lowest nesting level, you can use the <code>last</code> keyword. This is a Postgres Pro extension of the SQL/JSON standard.
<code>[subscript, ...]</code>	Array element accessor. <i>subscript</i> can be given in two forms: <i>index</i> or <i>start_index to end_index</i> . The first form returns a single array element by its index. The second form returns an array slice by the range of indexes, including the elements that correspond to the provided <i>start_index</i> and <i>end_index</i> . The specified <i>index</i> can be an integer, as well as an expression returning a single numeric value, which is automatically cast to integer. Index zero corresponds to the first array element. You can also use the <code>last</code> keyword to denote the last array element, which is useful for handling arrays of unknown length.
<code>[*]</code>	Wildcard array element accessor that returns all array elements.

8.15. Arrays

Postgres Pro allows columns of a table to be defined as variable-length multidimensional arrays. Arrays of any built-in or user-defined base type, enum type, composite type, range type, or domain can be created.

8.15.1. Declaration of Array Types

To illustrate the use of array types, we create this table:

```
CREATE TABLE sal_emp (
    name          text,
    pay_by_quarter integer[],
    schedule      text[][]
);
```

As shown, an array data type is named by appending square brackets (`[]`) to the data type name of the array elements. The above command will create a table named `sal_emp` with a column of type `text` (`name`), a one-dimensional array of type `integer` (`pay_by_quarter`), which represents the employee's salary by quarter, and a two-dimensional array of `text` (`schedule`), which represents the employee's weekly schedule.

The syntax for `CREATE TABLE` allows the exact size of arrays to be specified, for example:

```
CREATE TABLE tictactoe (
    squares      integer[3][3]
);
```

However, the current implementation ignores any supplied array size limits, i.e., the behavior is the same as for arrays of unspecified length.

The current implementation does not enforce the declared number of dimensions either. Arrays of a particular element type are all considered to be of the same type, regardless of size or number of dimen-

sions. So, declaring the array size or number of dimensions in `CREATE TABLE` is simply documentation; it does not affect run-time behavior.

An alternative syntax, which conforms to the SQL standard by using the keyword `ARRAY`, can be used for one-dimensional arrays. `pay_by_quarter` could have been defined as:

```
pay_by_quarter integer ARRAY[4],
```

Or, if no array size is to be specified:

```
pay_by_quarter integer ARRAY,
```

As before, however, Postgres Pro does not enforce the size restriction in any case.

8.15.2. Array Value Input

To write an array value as a literal constant, enclose the element values within curly braces and separate them by commas. (If you know C, this is not unlike the C syntax for initializing structures.) You can put double quotes around any element value, and must do so if it contains commas or curly braces. (More details appear below.) Thus, the general format of an array constant is the following:

```
'{ val1 delim val2 delim ... }'
```

where *delim* is the delimiter character for the type, as recorded in its `pg_type` entry. Among the standard data types provided in the Postgres Pro distribution, all use a comma (,), except for type `box` which uses a semicolon (;). Each *val* is either a constant of the array element type, or a subarray. An example of an array constant is:

```
'{{1,2,3},{4,5,6},{7,8,9}}'
```

This constant is a two-dimensional, 3-by-3 array consisting of three subarrays of integers.

To set an element of an array constant to `NULL`, write `NULL` for the element value. (Any upper- or lower-case variant of `NULL` will do.) If you want an actual string value “`NULL`”, you must put double quotes around it.

(These kinds of array constants are actually only a special case of the generic type constants discussed in [Section 4.1.2.7](#). The constant is initially treated as a string and passed to the array input conversion routine. An explicit type specification might be necessary.)

Now we can show some `INSERT` statements:

```
INSERT INTO sal_emp
VALUES ('Bill',
       '{10000, 10000, 10000, 10000}',
       '{{"meeting", "lunch"}, {"training", "presentation"}}');
```

```
INSERT INTO sal_emp
VALUES ('Carol',
       '{20000, 25000, 25000, 25000}',
       '{{"breakfast", "consulting"}, {"meeting", "lunch"}}');
```

The result of the previous two inserts looks like this:

```
SELECT * FROM sal_emp;
 name |      pay_by_quarter      |      schedule
-----+-----+-----
 Bill | {10000,10000,10000,10000} | {{meeting,lunch},{training,presentation}}
 Carol | {20000,25000,25000,25000} | {{breakfast,consulting},{meeting,lunch}}
(2 rows)
```

Multidimensional arrays must have matching extents for each dimension. A mismatch causes an error, for example:

```
INSERT INTO sal_emp
```

```
VALUES ('Bill',
        '{10000, 10000, 10000, 10000}',
        '{{"meeting", "lunch"}, {"meeting"}}');
ERROR:  malformed array literal: "{{"meeting", "lunch"}, {"meeting"}}"
DETAIL:  Multidimensional arrays must have sub-arrays with matching dimensions.
```

The `ARRAY` constructor syntax can also be used:

```
INSERT INTO sal_emp
VALUES ('Bill',
        ARRAY[10000, 10000, 10000, 10000],
        ARRAY[['meeting', 'lunch'], ['training', 'presentation']]);

INSERT INTO sal_emp
VALUES ('Carol',
        ARRAY[20000, 25000, 25000, 25000],
        ARRAY[['breakfast', 'consulting'], ['meeting', 'lunch']]);
```

Notice that the array elements are ordinary SQL constants or expressions; for instance, string literals are single quoted, instead of double quoted as they would be in an array literal. The `ARRAY` constructor syntax is discussed in more detail in [Section 4.2.12](#).

8.15.3. Accessing Arrays

Now, we can run some queries on the table. First, we show how to access a single element of an array. This query retrieves the names of the employees whose pay changed in the second quarter:

```
SELECT name FROM sal_emp WHERE pay_by_quarter[1] <> pay_by_quarter[2];
```

```
name
-----
Carol
(1 row)
```

The array subscript numbers are written within square brackets. By default Postgres Pro uses a one-based numbering convention for arrays, that is, an array of n elements starts with `array[1]` and ends with `array[n]`.

This query retrieves the third quarter pay of all employees:

```
SELECT pay_by_quarter[3] FROM sal_emp;
```

```
pay_by_quarter
-----
10000
25000
(2 rows)
```

We can also access arbitrary rectangular slices of an array, or subarrays. An array slice is denoted by writing *lower-bound:upper-bound* for one or more array dimensions. For example, this query retrieves the first item on Bill's schedule for the first two days of the week:

```
SELECT schedule[1:2][1:1] FROM sal_emp WHERE name = 'Bill';
```

```
schedule
-----
{{meeting},{training}}
(1 row)
```

If any dimension is written as a slice, i.e., contains a colon, then all dimensions are treated as slices. Any dimension that has only a single number (no colon) is treated as being from 1 to the number specified. For example, `[2]` is treated as `[1:2]`, as in this example:


```
SELECT schedule[1:2][2] FROM sal_emp WHERE name = 'Bill';
```

```
      schedule
-----
{{meeting,lunch},{training,presentation}}
(1 row)
```

To avoid confusion with the non-slice case, it's best to use slice syntax for all dimensions, e.g., [1:2][1:1], not [2][1:1].

It is possible to omit the *lower-bound* and/or *upper-bound* of a slice specifier; the missing bound is replaced by the lower or upper limit of the array's subscripts. For example:

```
SELECT schedule[:2][2:] FROM sal_emp WHERE name = 'Bill';
```

```
      schedule
-----
{{lunch},{presentation}}
(1 row)
```

```
SELECT schedule[:,1:1] FROM sal_emp WHERE name = 'Bill';
```

```
      schedule
-----
{{meeting},{training}}
(1 row)
```

An array subscript expression will return null if either the array itself or any of the subscript expressions are null. Also, null is returned if a subscript is outside the array bounds (this case does not raise an error). For example, if `schedule` currently has the dimensions [1:3][1:2] then referencing `schedule[3][3]` yields NULL. Similarly, an array reference with the wrong number of subscripts yields a null rather than an error.

An array slice expression likewise yields null if the array itself or any of the subscript expressions are null. However, in other cases such as selecting an array slice that is completely outside the current array bounds, a slice expression yields an empty (zero-dimensional) array instead of null. (This does not match non-slice behavior and is done for historical reasons.) If the requested slice partially overlaps the array bounds, then it is silently reduced to just the overlapping region instead of returning null.

The current dimensions of any array value can be retrieved with the `array_dims` function:

```
SELECT array_dims(schedule) FROM sal_emp WHERE name = 'Carol';
```

```
      array_dims
-----
[1:2][1:2]
(1 row)
```

`array_dims` produces a text result, which is convenient for people to read but perhaps inconvenient for programs. Dimensions can also be retrieved with `array_upper` and `array_lower`, which return the upper and lower bound of a specified array dimension, respectively:

```
SELECT array_upper(schedule, 1) FROM sal_emp WHERE name = 'Carol';
```

```
      array_upper
-----
2
(1 row)
```

`array_length` will return the length of a specified array dimension:

```
SELECT array_length(schedule, 1) FROM sal_emp WHERE name = 'Carol';
```

```
array_length
-----
                2
(1 row)
```

`cardinality` returns the total number of elements in an array across all dimensions. It is effectively the number of rows a call to `unnest` would yield:

```
SELECT cardinality(schedule) FROM sal_emp WHERE name = 'Carol';

cardinality
-----
                4
(1 row)
```

8.15.4. Modifying Arrays

An array value can be replaced completely:

```
UPDATE sal_emp SET pay_by_quarter = '{25000,25000,27000,27000}'
WHERE name = 'Carol';
```

or using the `ARRAY` expression syntax:

```
UPDATE sal_emp SET pay_by_quarter = ARRAY[25000,25000,27000,27000]
WHERE name = 'Carol';
```

An array can also be updated at a single element:

```
UPDATE sal_emp SET pay_by_quarter[4] = 15000
WHERE name = 'Bill';
```

or updated in a slice:

```
UPDATE sal_emp SET pay_by_quarter[1:2] = '{27000,27000}'
WHERE name = 'Carol';
```

The slice syntaxes with omitted *lower-bound* and/or *upper-bound* can be used too, but only when updating an array value that is not `NULL` or zero-dimensional (otherwise, there is no existing subscript limit to substitute).

A stored array value can be enlarged by assigning to elements not already present. Any positions between those previously present and the newly assigned elements will be filled with nulls. For example, if array `myarray` currently has 4 elements, it will have six elements after an update that assigns to `myarray[6]`; `myarray[5]` will contain null. Currently, enlargement in this fashion is only allowed for one-dimensional arrays, not multidimensional arrays.

Subscripted assignment allows creation of arrays that do not use one-based subscripts. For example one might assign to `myarray[-2:7]` to create an array with subscript values from -2 to 7.

New array values can also be constructed using the concatenation operator, `||`:

```
SELECT ARRAY[1,2] || ARRAY[3,4];
?column?
-----
{1,2,3,4}
(1 row)

SELECT ARRAY[5,6] || ARRAY[[1,2],[3,4]];
?column?
-----
{{5,6},{1,2},{3,4}}
```

```
(1 row)
```

The concatenation operator allows a single element to be pushed onto the beginning or end of a one-dimensional array. It also accepts two N -dimensional arrays, or an N -dimensional and an $N+1$ -dimensional array.

When a single element is pushed onto either the beginning or end of a one-dimensional array, the result is an array with the same lower bound subscript as the array operand. For example:

```
SELECT array_dims(1 || '[0:1]={2,3}'::int[]);
 array_dims
-----
 [0:2]
(1 row)
```

```
SELECT array_dims(ARRAY[1,2] || 3);
 array_dims
-----
 [1:3]
(1 row)
```

When two arrays with an equal number of dimensions are concatenated, the result retains the lower bound subscript of the left-hand operand's outer dimension. The result is an array comprising every element of the left-hand operand followed by every element of the right-hand operand. For example:

```
SELECT array_dims(ARRAY[1,2] || ARRAY[3,4,5]);
 array_dims
-----
 [1:5]
(1 row)

SELECT array_dims(ARRAY[[1,2],[3,4]] || ARRAY[[5,6],[7,8],[9,0]]);
 array_dims
-----
 [1:5][1:2]
(1 row)
```

When an N -dimensional array is pushed onto the beginning or end of an $N+1$ -dimensional array, the result is analogous to the element-array case above. Each N -dimensional sub-array is essentially an element of the $N+1$ -dimensional array's outer dimension. For example:

```
SELECT array_dims(ARRAY[1,2] || ARRAY[[3,4],[5,6]]);
 array_dims
-----
 [1:3][1:2]
(1 row)
```

An array can also be constructed by using the functions `array_prepend`, `array_append`, or `array_cat`. The first two only support one-dimensional arrays, but `array_cat` supports multidimensional arrays. Some examples:

```
SELECT array_prepend(1, ARRAY[2,3]);
 array_prepend
-----
 {1,2,3}
(1 row)

SELECT array_append(ARRAY[1,2], 3);
 array_append
-----
 {1,2,3}
```

```
(1 row)
```

```
SELECT array_cat (ARRAY[1,2], ARRAY[3,4]);
array_cat
-----
{1,2,3,4}
(1 row)
```

```
SELECT array_cat (ARRAY[[1,2],[3,4]], ARRAY[5,6]);
array_cat
-----
{{1,2},{3,4},{5,6}}
(1 row)
```

```
SELECT array_cat (ARRAY[5,6], ARRAY[[1,2],[3,4]]);
array_cat
-----
{{5,6},{1,2},{3,4}}
```

In simple cases, the concatenation operator discussed above is preferred over direct use of these functions. However, because the concatenation operator is overloaded to serve all three cases, there are situations where use of one of the functions is helpful to avoid ambiguity. For example consider:

```
SELECT ARRAY[1, 2] || '{3, 4}'; -- the untyped literal is taken as an array
?column?
-----
{1,2,3,4}
```

```
SELECT ARRAY[1, 2] || '7'; -- so is this one
ERROR:  malformed array literal: "7"
```

```
SELECT ARRAY[1, 2] || NULL; -- so is an undecorated NULL
?column?
-----
{1,2}
(1 row)
```

```
SELECT array_append(ARRAY[1, 2], NULL); -- this might have been meant
array_append
-----
{1,2,NULL}
```

In the examples above, the parser sees an integer array on one side of the concatenation operator, and a constant of undetermined type on the other. The heuristic it uses to resolve the constant's type is to assume it's of the same type as the operator's other input — in this case, integer array. So the concatenation operator is presumed to represent `array_cat`, not `array_append`. When that's the wrong choice, it could be fixed by casting the constant to the array's element type; but explicit use of `array_append` might be a preferable solution.

8.15.5. Searching in Arrays

To search for a value in an array, each value must be checked. This can be done manually, if you know the size of the array. For example:

```
SELECT * FROM sal_emp WHERE pay_by_quarter[1] = 10000 OR
                             pay_by_quarter[2] = 10000 OR
                             pay_by_quarter[3] = 10000 OR
                             pay_by_quarter[4] = 10000;
```

However, this quickly becomes tedious for large arrays, and is not helpful if the size of the array is unknown. An alternative method is described in [Section 9.25](#). The above query could be replaced by:

```
SELECT * FROM sal_emp WHERE 10000 = ANY (pay_by_quarter);
```

In addition, you can find rows where the array has all values equal to 10000 with:

```
SELECT * FROM sal_emp WHERE 10000 = ALL (pay_by_quarter);
```

Alternatively, the `generate_subscripts` function can be used. For example:

```
SELECT * FROM
  (SELECT pay_by_quarter,
    generate_subscripts(pay_by_quarter, 1) AS s
   FROM sal_emp) AS foo
WHERE pay_by_quarter[s] = 10000;
```

This function is described in [Table 9.68](#).

You can also search an array using the `&&` operator, which checks whether the left operand overlaps with the right operand. For instance:

```
SELECT * FROM sal_emp WHERE pay_by_quarter && ARRAY[10000];
```

This and other array operators are further described in [Section 9.19](#). It can be accelerated by an appropriate index, as described in [Section 11.2](#).

You can also search for specific values in an array using the `array_position` and `array_positions` functions. The former returns the subscript of the first occurrence of a value in an array; the latter returns an array with the subscripts of all occurrences of the value in the array. For example:

```
SELECT array_position(ARRAY['sun','mon','tue','wed','thu','fri','sat'], 'mon');
array_position
-----
                2
(1 row)

SELECT array_positions(ARRAY[1, 4, 3, 1, 3, 4, 2, 1], 1);
array_positions
-----
{1,4,8}
(1 row)
```

Tip

Arrays are not sets; searching for specific array elements can be a sign of database misdesign. Consider using a separate table with a row for each item that would be an array element. This will be easier to search, and is likely to scale better for a large number of elements.

8.15.6. Array Input and Output Syntax

The external text representation of an array value consists of items that are interpreted according to the I/O conversion rules for the array's element type, plus decoration that indicates the array structure. The decoration consists of curly braces (`{` and `}`) around the array value plus delimiter characters between adjacent items. The delimiter character is usually a comma (`,`) but can be something else: it is determined by the `typdelim` setting for the array's element type. Among the standard data types provided in the Postgres Pro distribution, all use a comma, except for type `box`, which uses a semicolon (`;`). In a multidimensional array, each dimension (row, plane, cube, etc.) gets its own level of curly braces, and delimiters must be written between adjacent curly-braced entities of the same level.

The array output routine will put double quotes around element values if they are empty strings, contain curly braces, delimiter characters, double quotes, backslashes, or white space, or match the word `NULL`. Double quotes and backslashes embedded in element values will be backslash-escaped. For numeric

data types it is safe to assume that double quotes will never appear, but for textual data types one should be prepared to cope with either the presence or absence of quotes.

By default, the lower bound index value of an array's dimensions is set to one. To represent arrays with other lower bounds, the array subscript ranges can be specified explicitly before writing the array contents. This decoration consists of square brackets (`[]`) around each array dimension's lower and upper bounds, with a colon (`:`) delimiter character in between. The array dimension decoration is followed by an equal sign (`=`). For example:

```
SELECT f1[1][-2][3] AS e1, f1[1][-1][5] AS e2
FROM (SELECT '[1:1][-2:-1][3:5]={{{1,2,3},{4,5,6}}}'::int[] AS f1) AS ss;
```

```
e1 | e2
----+----
 1 |  6
(1 row)
```

The array output routine will include explicit dimensions in its result only when there are one or more lower bounds different from one.

If the value written for an element is `NULL` (in any case variant), the element is taken to be `NULL`. The presence of any quotes or backslashes disables this and allows the literal string value “`NULL`” to be entered. Also, for backward compatibility with pre-8.2 versions of PostgreSQL, the [array_nulls](#) configuration parameter can be turned `off` to suppress recognition of `NULL` as a `NULL`.

As shown previously, when writing an array value you can use double quotes around any individual array element. You *must* do so if the element value would otherwise confuse the array-value parser. For example, elements containing curly braces, commas (or the data type's delimiter character), double quotes, backslashes, or leading or trailing whitespace must be double-quoted. Empty strings and strings matching the word `NULL` must be quoted, too. To put a double quote or backslash in a quoted array element value, precede it with a backslash. Alternatively, you can avoid quotes and use backslash-escaping to protect all data characters that would otherwise be taken as array syntax.

You can add whitespace before a left brace or after a right brace. You can also add whitespace before or after any individual item string. In all of these cases the whitespace will be ignored. However, whitespace within double-quoted elements, or surrounded on both sides by non-whitespace characters of an element, is not ignored.

Tip

The `ARRAY` constructor syntax (see [Section 4.2.12](#)) is often easier to work with than the array-literal syntax when writing array values in SQL commands. In `ARRAY`, individual element values are written the same way they would be written when not members of an array.

8.16. Composite Types

A *composite type* represents the structure of a row or record; it is essentially just a list of field names and their data types. Postgres Pro allows composite types to be used in many of the same ways that simple types can be used. For example, a column of a table can be declared to be of a composite type.

8.16.1. Declaration of Composite Types

Here are two simple examples of defining composite types:

```
CREATE TYPE complex AS (
    r      double precision,
    i      double precision
);
```

```
CREATE TYPE inventory_item AS (  
    name            text,  
    supplier_id     integer,  
    price           numeric  
);
```

The syntax is comparable to `CREATE TABLE`, except that only field names and types can be specified; no constraints (such as `NOT NULL`) can presently be included. Note that the `AS` keyword is essential; without it, the system will think a different kind of `CREATE TYPE` command is meant, and you will get odd syntax errors.

Having defined the types, we can use them to create tables:

```
CREATE TABLE on_hand (  
    item            inventory_item,  
    count          integer  
);
```

```
INSERT INTO on_hand VALUES (ROW('fuzzy dice', 42, 1.99), 1000);
```

or functions:

```
CREATE FUNCTION price_extension(inventory_item, integer) RETURNS numeric  
AS 'SELECT $1.price * $2' LANGUAGE SQL;
```

```
SELECT price_extension(item, 10) FROM on_hand;
```

Whenever you create a table, a composite type is also automatically created, with the same name as the table, to represent the table's row type. For example, had we said:

```
CREATE TABLE inventory_item (  
    name            text,  
    supplier_id     integer REFERENCES suppliers,  
    price           numeric CHECK (price > 0)  
);
```

then the same `inventory_item` composite type shown above would come into being as a byproduct, and could be used just as above. Note however an important restriction of the current implementation: since no constraints are associated with a composite type, the constraints shown in the table definition *do not apply* to values of the composite type outside the table. (To work around this, create a [domain](#) over the composite type, and apply the desired constraints as `CHECK` constraints of the domain.)

8.16.2. Constructing Composite Values

To write a composite value as a literal constant, enclose the field values within parentheses and separate them by commas. You can put double quotes around any field value, and must do so if it contains commas or parentheses. (More details appear [below](#).) Thus, the general format of a composite constant is the following:

```
'( val1 , val2 , ... )'
```

An example is:

```
'("fuzzy dice",42,1.99)'
```

which would be a valid value of the `inventory_item` type defined above. To make a field be `NULL`, write no characters at all in its position in the list. For example, this constant specifies a `NULL` third field:

```
'("fuzzy dice",42,)'
```

If you want an empty string rather than `NULL`, write double quotes:

```
'("",42,)'
```

Here the first field is a non-NULL empty string, the third is NULL.

(These constants are actually only a special case of the generic type constants discussed in [Section 4.1.2.7](#). The constant is initially treated as a string and passed to the composite-type input conversion routine. An explicit type specification might be necessary to tell which type to convert the constant to.)

The `ROW` expression syntax can also be used to construct composite values. In most cases this is considerably simpler to use than the string-literal syntax since you don't have to worry about multiple layers of quoting. We already used this method above:

```
ROW('fuzzy dice', 42, 1.99)
ROW('', 42, NULL)
```

The `ROW` keyword is actually optional as long as you have more than one field in the expression, so these can be simplified to:

```
('fuzzy dice', 42, 1.99)
('', 42, NULL)
```

The `ROW` expression syntax is discussed in more detail in [Section 4.2.13](#).

8.16.3. Accessing Composite Types

To access a field of a composite column, one writes a dot and the field name, much like selecting a field from a table name. In fact, it's so much like selecting from a table name that you often have to use parentheses to keep from confusing the parser. For example, you might try to select some subfields from our `on_hand` example table with something like:

```
SELECT item.name FROM on_hand WHERE item.price > 9.99;
```

This will not work since the name `item` is taken to be a table name, not a column name of `on_hand`, per SQL syntax rules. You must write it like this:

```
SELECT (item).name FROM on_hand WHERE (item).price > 9.99;
```

or if you need to use the table name as well (for instance in a multitable query), like this:

```
SELECT (on_hand.item).name FROM on_hand WHERE (on_hand.item).price > 9.99;
```

Now the parenthesized object is correctly interpreted as a reference to the `item` column, and then the subfield can be selected from it.

Similar syntactic issues apply whenever you select a field from a composite value. For instance, to select just one field from the result of a function that returns a composite value, you'd need to write something like:

```
SELECT (my_func(...)).field FROM ...
```

Without the extra parentheses, this will generate a syntax error.

The special field name `*` means “all fields”, as further explained in [Section 8.16.5](#).

8.16.4. Modifying Composite Types

Here are some examples of the proper syntax for inserting and updating composite columns. First, inserting or updating a whole column:

```
INSERT INTO mytab (complex_col) VALUES((1.1,2.2));

UPDATE mytab SET complex_col = ROW(1.1,2.2) WHERE ...;
```

The first example omits `ROW`, the second uses it; we could have done it either way.

We can update an individual subfield of a composite column:

```
UPDATE mytab SET complex_col.r = (complex_col).r + 1 WHERE ...;
```


Notice here that we don't need to (and indeed cannot) put parentheses around the column name appearing just after `SET`, but we do need parentheses when referencing the same column in the expression to the right of the equal sign.

And we can specify subfields as targets for `INSERT`, too:

```
INSERT INTO mytab (complex_col.r, complex_col.i) VALUES(1.1, 2.2);
```

Had we not supplied values for all the subfields of the column, the remaining subfields would have been filled with null values.

8.16.5. Using Composite Types in Queries

There are various special syntax rules and behaviors associated with composite types in queries. These rules provide useful shortcuts, but can be confusing if you don't know the logic behind them.

In Postgres Pro, a reference to a table name (or alias) in a query is effectively a reference to the composite value of the table's current row. For example, if we had a table `inventory_item` as shown [above](#), we could write:

```
SELECT c FROM inventory_item c;
```

This query produces a single composite-valued column, so we might get output like:

```
      c
-----
("fuzzy dice",42,1.99)
(1 row)
```

Note however that simple names are matched to column names before table names, so this example works only because there is no column named `c` in the query's tables.

The ordinary qualified-column-name syntax `table_name.column_name` can be understood as applying [field selection](#) to the composite value of the table's current row. (For efficiency reasons, it's not actually implemented that way.)

When we write

```
SELECT c.* FROM inventory_item c;
```

then, according to the SQL standard, we should get the contents of the table expanded into separate columns:

```
   name   | supplier_id | price
-----+-----+-----
fuzzy dice |           42 |   1.99
(1 row)
```

as if the query were

```
SELECT c.name, c.supplier_id, c.price FROM inventory_item c;
```

Postgres Pro will apply this expansion behavior to any composite-valued expression, although as shown [above](#), you need to write parentheses around the value that `.*` is applied to whenever it's not a simple table name. For example, if `myfunc()` is a function returning a composite type with columns `a`, `b`, and `c`, then these two queries have the same result:

```
SELECT (myfunc(x)).* FROM some_table;
SELECT (myfunc(x)).a, (myfunc(x)).b, (myfunc(x)).c FROM some_table;
```

Tip

Postgres Pro handles column expansion by actually transforming the first form into the second. So, in this example, `myfunc()` would get invoked three times per row with either syntax. If it's an expensive function you may wish to avoid that, which you can do with a query like:

```
SELECT m.* FROM some_table, LATERAL myfunc(x) AS m;
```

Placing the function in a `LATERAL FROM` item keeps it from being invoked more than once per row. `m.*` is still expanded into `m.a`, `m.b`, `m.c`, but now those variables are just references to the output of the `FROM` item. (The `LATERAL` keyword is optional here, but we show it to clarify that the function is getting `x` from `some_table`.)

The `composite_value.*` syntax results in column expansion of this kind when it appears at the top level of a [SELECT output list](#), a [RETURNING list](#) in `INSERT/UPDATE/DELETE/MERGE`, a [VALUES clause](#), or a [row constructor](#). In all other contexts (including when nested inside one of those constructs), attaching `.*` to a composite value does not change the value, since it means “all columns” and so the same composite value is produced again. For example, if `somefunc()` accepts a composite-valued argument, these queries are the same:

```
SELECT somefunc(c.*) FROM inventory_item c;
SELECT somefunc(c) FROM inventory_item c;
```

In both cases, the current row of `inventory_item` is passed to the function as a single composite-valued argument. Even though `.*` does nothing in such cases, using it is good style, since it makes clear that a composite value is intended. In particular, the parser will consider `c` in `c.*` to refer to a table name or alias, not to a column name, so that there is no ambiguity; whereas without `.*`, it is not clear whether `c` means a table name or a column name, and in fact the column-name interpretation will be preferred if there is a column named `c`.

Another example demonstrating these concepts is that all these queries mean the same thing:

```
SELECT * FROM inventory_item c ORDER BY c;
SELECT * FROM inventory_item c ORDER BY c.*;
SELECT * FROM inventory_item c ORDER BY ROW(c.*);
```

All of these `ORDER BY` clauses specify the row's composite value, resulting in sorting the rows according to the rules described in [Section 9.25.6](#). However, if `inventory_item` contained a column named `c`, the first case would be different from the others, as it would mean to sort by that column only. Given the column names previously shown, these queries are also equivalent to those above:

```
SELECT * FROM inventory_item c ORDER BY ROW(c.name, c.supplier_id, c.price);
SELECT * FROM inventory_item c ORDER BY (c.name, c.supplier_id, c.price);
```

(The last case uses a row constructor with the key word `ROW` omitted.)

Another special syntactical behavior associated with composite values is that we can use *functional notation* for extracting a field of a composite value. The simple way to explain this is that the notations `field(table)` and `table.field` are interchangeable. For example, these queries are equivalent:

```
SELECT c.name FROM inventory_item c WHERE c.price > 1000;
SELECT name(c) FROM inventory_item c WHERE price(c) > 1000;
```

Moreover, if we have a function that accepts a single argument of a composite type, we can call it with either notation. These queries are all equivalent:

```
SELECT somefunc(c) FROM inventory_item c;
SELECT somefunc(c.*) FROM inventory_item c;
SELECT c.somefunc FROM inventory_item c;
```

This equivalence between functional notation and field notation makes it possible to use functions on composite types to implement “computed fields”. An application using the last query above wouldn't need to be directly aware that `somefunc` isn't a real column of the table.

Tip

Because of this behavior, it's unwise to give a function that takes a single composite-type argument the same name as any of the fields of that composite type. If there is ambiguity, the field-name

interpretation will be chosen if field-name syntax is used, while the function will be chosen if function-call syntax is used. However, Postgres Pro versions before 11 always chose the field-name interpretation, unless the syntax of the call required it to be a function call. One way to force the function interpretation in older versions is to schema-qualify the function name, that is, write `schema.func(compositevalue)`.

8.16.6. Composite Type Input and Output Syntax

The external text representation of a composite value consists of items that are interpreted according to the I/O conversion rules for the individual field types, plus decoration that indicates the composite structure. The decoration consists of parentheses ((and)) around the whole value, plus commas (,) between adjacent items. Whitespace outside the parentheses is ignored, but within the parentheses it is considered part of the field value, and might or might not be significant depending on the input conversion rules for the field data type. For example, in:

```
' ( 42) '
```

the whitespace will be ignored if the field type is integer, but not if it is text.

As shown previously, when writing a composite value you can write double quotes around any individual field value. You *must* do so if the field value would otherwise confuse the composite-value parser. In particular, fields containing parentheses, commas, double quotes, or backslashes must be double-quoted. To put a double quote or backslash in a quoted composite field value, precede it with a backslash. (Also, a pair of double quotes within a double-quoted field value is taken to represent a double quote character, analogously to the rules for single quotes in SQL literal strings.) Alternatively, you can avoid quoting and use backslash-escaping to protect all data characters that would otherwise be taken as composite syntax.

A completely empty field value (no characters at all between the commas or parentheses) represents a NULL. To write a value that is an empty string rather than NULL, write "".

The composite output routine will put double quotes around field values if they are empty strings or contain parentheses, commas, double quotes, backslashes, or white space. (Doing so for white space is not essential, but aids legibility.) Double quotes and backslashes embedded in field values will be doubled.

Note

Remember that what you write in an SQL command will first be interpreted as a string literal, and then as a composite. This doubles the number of backslashes you need (assuming escape string syntax is used). For example, to insert a `text` field containing a double quote and a backslash in a composite value, you'd need to write:

```
INSERT ... VALUES ('("\\""\")');
```

The string-literal processor removes one level of backslashes, so that what arrives at the composite-value parser looks like ("\""). In turn, the string fed to the `text` data type's input routine becomes "\. (If we were working with a data type whose input routine also treated backslashes specially, `bytea` for example, we might need as many as eight backslashes in the command to get one backslash into the stored composite field.) Dollar quoting (see [Section 4.1.2.4](#)) can be used to avoid the need to double backslashes.

Tip

The `ROW` constructor syntax is usually easier to work with than the composite-literal syntax when writing composite values in SQL commands. In `ROW`, individual field values are written the same way they would be written when not members of a composite.

8.17. Range Types

Range types are data types representing a range of values of some element type (called the range's *subtype*). For instance, ranges of `timestamp` might be used to represent the ranges of time that a meeting room is reserved. In this case the data type is `tsrange` (short for “timestamp range”), and `timestamp` is the subtype. The subtype must have a total order so that it is well-defined whether element values are within, before, or after a range of values.

Range types are useful because they represent many element values in a single range value, and because concepts such as overlapping ranges can be expressed clearly. The use of time and date ranges for scheduling purposes is the clearest example; but price ranges, measurement ranges from an instrument, and so forth can also be useful.

Every range type has a corresponding multirange type. A multirange is an ordered list of non-contiguous, non-empty, non-null ranges. Most range operators also work on multiranges, and they have a few functions of their own.

8.17.1. Built-in Range and Multirange Types

Postgres Pro comes with the following built-in range types:

- `int4range` — Range of integer, `int4multirange` — corresponding Multirange
- `int8range` — Range of bigint, `int8multirange` — corresponding Multirange
- `numrange` — Range of numeric, `nummultirange` — corresponding Multirange
- `tsrange` — Range of timestamp without time zone, `tsmultirange` — corresponding Multirange
- `tstzrange` — Range of timestamp with time zone, `tstzmultirange` — corresponding Multirange
- `daterange` — Range of date, `datemultirange` — corresponding Multirange

In addition, you can define your own range types; see [CREATE TYPE](#) for more information.

8.17.2. Examples

```
CREATE TABLE reservation (room int, during tsrange);
INSERT INTO reservation VALUES
    (1108, '[2010-01-01 14:30, 2010-01-01 15:30)');

-- Containment
SELECT int4range(10, 20) @> 3;

-- Overlaps
SELECT numrange(11.1, 22.2) && numrange(20.0, 30.0);

-- Extract the upper bound
SELECT upper(int8range(15, 25));

-- Compute the intersection
SELECT int4range(10, 20) * int4range(15, 25);

-- Is the range empty?
SELECT isempty(numrange(1, 5));
```

See [Table 9.56](#) and [Table 9.58](#) for complete lists of operators and functions on range types.

8.17.3. Inclusive and Exclusive Bounds

Every non-empty range has two bounds, the lower bound and the upper bound. All points between these values are included in the range. An inclusive bound means that the boundary point itself is included in the range as well, while an exclusive bound means that the boundary point is not included in the range.

In the text form of a range, an inclusive lower bound is represented by “[” while an exclusive lower bound is represented by “(”. Likewise, an inclusive upper bound is represented by “]”, while an exclusive upper bound is represented by “)”. (See [Section 8.17.5](#) for more details.)

The functions `lower_inc` and `upper_inc` test the inclusivity of the lower and upper bounds of a range value, respectively.

8.17.4. Infinite (Unbounded) Ranges

The lower bound of a range can be omitted, meaning that all values less than the upper bound are included in the range, e.g., `(, 3]`. Likewise, if the upper bound of the range is omitted, then all values greater than the lower bound are included in the range. If both lower and upper bounds are omitted, all values of the element type are considered to be in the range. Specifying a missing bound as inclusive is automatically converted to exclusive, e.g., `[,]` is converted to `(,)`. You can think of these missing values as +/-infinity, but they are special range type values and are considered to be beyond any range element type's +/-infinity values.

Element types that have the notion of “infinity” can use them as explicit bound values. For example, with timestamp ranges, `[today, infinity)` excludes the special timestamp value `infinity`, while `[today, infinity]` include it, as does `[today,)` and `[today,]`.

The functions `lower_inf` and `upper_inf` test for infinite lower and upper bounds of a range, respectively.

8.17.5. Range Input/Output

The input for a range value must follow one of the following patterns:

```
(lower-bound, upper-bound)
(lower-bound, upper-bound]
[lower-bound, upper-bound)
[lower-bound, upper-bound]
empty
```

The parentheses or brackets indicate whether the lower and upper bounds are exclusive or inclusive, as described previously. Notice that the final pattern is `empty`, which represents an empty range (a range that contains no points).

The *lower-bound* may be either a string that is valid input for the subtype, or empty to indicate no lower bound. Likewise, *upper-bound* may be either a string that is valid input for the subtype, or empty to indicate no upper bound.

Each bound value can be quoted using " (double quote) characters. This is necessary if the bound value contains parentheses, brackets, commas, double quotes, or backslashes, since these characters would otherwise be taken as part of the range syntax. To put a double quote or backslash in a quoted bound value, precede it with a backslash. (Also, a pair of double quotes within a double-quoted bound value is taken to represent a double quote character, analogously to the rules for single quotes in SQL literal strings.) Alternatively, you can avoid quoting and use backslash-escaping to protect all data characters that would otherwise be taken as range syntax. Also, to write a bound value that is an empty string, write "", since writing nothing means an infinite bound.

Whitespace is allowed before and after the range value, but any whitespace between the parentheses or brackets is taken as part of the lower or upper bound value. (Depending on the element type, it might or might not be significant.)

Note

These rules are very similar to those for writing field values in composite-type literals. See [Section 8.16.6](#) for additional commentary.

Examples:

```
-- includes 3, does not include 7, and does include all points in between
SELECT '[3,7) '::int4range;

-- does not include either 3 or 7, but includes all points in between
SELECT '(3,7) '::int4range;

-- includes only the single point 4
SELECT '[4,4] '::int4range;

-- includes no points (and will be normalized to 'empty')
SELECT '[4,4) '::int4range;
```

The input for a multirange is curly brackets (`{` and `}`) containing zero or more valid ranges, separated by commas. Whitespace is permitted around the brackets and commas. This is intended to be reminiscent of array syntax, although multiranges are much simpler: they have just one dimension and there is no need to quote their contents. (The bounds of their ranges may be quoted as above however.)

Examples:

```
SELECT '{} '::int4multirange;
SELECT '{[3,7)} '::int4multirange;
SELECT '{[3,7), [8,9)} '::int4multirange;
```

8.17.6. Constructing Ranges and Multiranges

Each range type has a constructor function with the same name as the range type. Using the constructor function is frequently more convenient than writing a range literal constant, since it avoids the need for extra quoting of the bound values. The constructor function accepts two or three arguments. The two-argument form constructs a range in standard form (lower bound inclusive, upper bound exclusive), while the three-argument form constructs a range with bounds of the form specified by the third argument. The third argument must be one of the strings `"()"`, `"[]"`, `"[]"`, or `"[]"`. For example:

```
-- The full form is: lower bound, upper bound, and text argument indicating
-- inclusivity/exclusivity of bounds.
SELECT numrange(1.0, 14.0, '()');

-- If the third argument is omitted, '()' is assumed.
SELECT numrange(1.0, 14.0);

-- Although '()' is specified here, on display the value will be converted to
-- canonical form, since int8range is a discrete range type (see below).
SELECT int8range(1, 14, '()');

-- Using NULL for either bound causes the range to be unbounded on that side.
SELECT numrange(NULL, 2.2);
```

Each range type also has a multirange constructor with the same name as the multirange type. The constructor function takes zero or more arguments which are all ranges of the appropriate type. For example:

```
SELECT nummultirange();
SELECT nummultirange(numrange(1.0, 14.0));
SELECT nummultirange(numrange(1.0, 14.0), numrange(20.0, 25.0));
```

8.17.7. Discrete Range Types

A discrete range is one whose element type has a well-defined “step”, such as `integer` or `date`. In these types two elements can be said to be adjacent, when there are no valid values between them. This contrasts with continuous ranges, where it's always (or almost always) possible to identify other element

values between two given values. For example, a range over the `numeric` type is continuous, as is a range over `timestamp`. (Even though `timestamp` has limited precision, and so could theoretically be treated as discrete, it's better to consider it continuous since the step size is normally not of interest.)

Another way to think about a discrete range type is that there is a clear idea of a “next” or “previous” value for each element value. Knowing that, it is possible to convert between inclusive and exclusive representations of a range's bounds, by choosing the next or previous element value instead of the one originally given. For example, in an integer range type `[4, 8]` and `(3, 9)` denote the same set of values; but this would not be so for a range over `numeric`.

A discrete range type should have a *canonicalization* function that is aware of the desired step size for the element type. The canonicalization function is charged with converting equivalent values of the range type to have identical representations, in particular consistently inclusive or exclusive bounds. If a canonicalization function is not specified, then ranges with different formatting will always be treated as unequal, even though they might represent the same set of values in reality.

The built-in range types `int4range`, `int8range`, and `daterange` all use a canonical form that includes the lower bound and excludes the upper bound; that is, `()`. User-defined range types can use other conventions, however.

8.17.8. Defining New Range Types

Users can define their own range types. The most common reason to do this is to use ranges over subtypes not provided among the built-in range types. For example, to define a new range type of subtype `float8`:

```
CREATE TYPE floatrange AS RANGE (  
    subtype = float8,  
    subtype_diff = float8mi  
);  
  
SELECT '[1.234, 5.678]':floatrange;
```

Because `float8` has no meaningful “step”, we do not define a canonicalization function in this example.

When you define your own range you automatically get a corresponding `multirange` type.

Defining your own range type also allows you to specify a different subtype B-tree operator class or collation to use, so as to change the sort ordering that determines which values fall into a given range.

If the subtype is considered to have discrete rather than continuous values, the `CREATE TYPE` command should specify a `canonical` function. The canonicalization function takes an input range value, and must return an equivalent range value that may have different bounds and formatting. The canonical output for two ranges that represent the same set of values, for example the integer ranges `[1, 7]` and `[1, 8)`, must be identical. It doesn't matter which representation you choose to be the canonical one, so long as two equivalent values with different formattings are always mapped to the same value with the same formatting. In addition to adjusting the inclusive/exclusive bounds format, a canonicalization function might round off boundary values, in case the desired step size is larger than what the subtype is capable of storing. For instance, a range type over `timestamp` could be defined to have a step size of an hour, in which case the canonicalization function would need to round off bounds that weren't a multiple of an hour, or perhaps throw an error instead.

In addition, any range type that is meant to be used with GiST or SP-GiST indexes should define a subtype difference, or `subtype_diff`, function. (The index will still work without `subtype_diff`, but it is likely to be considerably less efficient than if a difference function is provided.) The subtype difference function takes two input values of the subtype, and returns their difference (i.e., x minus y) represented as a `float8` value. In our example above, the function `float8mi` that underlies the regular `float8` minus operator can be used; but for any other subtype, some type conversion would be necessary. Some creative thought about how to represent differences as numbers might be needed, too. To the greatest extent possible, the `subtype_diff` function should agree with the sort ordering implied by the selected

operator class and collation; that is, its result should be positive whenever its first argument is greater than its second according to the sort ordering.

A less-oversimplified example of a `subtype_diff` function is:

```
CREATE FUNCTION time_subtype_diff(x time, y time) RETURNS float8 AS
'SELECT EXTRACT(EPOCH FROM (x - y))' LANGUAGE sql STRICT IMMUTABLE;
```

```
CREATE TYPE timerange AS RANGE (
    subtype = time,
    subtype_diff = time_subtype_diff
);
```

```
SELECT '[11:10, 23:00]':::timerange;
```

See [CREATE TYPE](#) for more information about creating range types.

8.17.9. Indexing

GiST and SP-GiST indexes can be created for table columns of range types. GiST indexes can be also created for table columns of multirange types. For instance, to create a GiST index:

```
CREATE INDEX reservation_idx ON reservation USING GIST (during);
```

A GiST or SP-GiST index on ranges can accelerate queries involving these range operators: `=`, `&&`, `<@`, `@>`, `<<`, `>>`, `-|-`, `&<`, and `&>`. A GiST index on multiranges can accelerate queries involving the same set of multirange operators. A GiST index on ranges and GiST index on multiranges can also accelerate queries involving these cross-type range to multirange and multirange to range operators correspondingly: `&&`, `<@`, `@>`, `<<`, `>>`, `-|-`, `&<`, and `&>`. See [Table 9.56](#) for more information.

In addition, B-tree and hash indexes can be created for table columns of range types. For these index types, basically the only useful range operation is equality. There is a B-tree sort ordering defined for range values, with corresponding `<` and `>` operators, but the ordering is rather arbitrary and not usually useful in the real world. Range types' B-tree and hash support is primarily meant to allow sorting and hashing internally in queries, rather than creation of actual indexes.

8.17.10. Constraints on Ranges

While `UNIQUE` is a natural constraint for scalar values, it is usually unsuitable for range types. Instead, an exclusion constraint is often more appropriate (see [CREATE TABLE ... CONSTRAINT ... EXCLUDE](#)). Exclusion constraints allow the specification of constraints such as “non-overlapping” on a range type. For example:

```
CREATE TABLE reservation (
    during tsrange,
    EXCLUDE USING GIST (during WITH &&)
);
```

That constraint will prevent any overlapping values from existing in the table at the same time:

```
INSERT INTO reservation VALUES
('[2010-01-01 11:30, 2010-01-01 15:00)');
INSERT 0 1

INSERT INTO reservation VALUES
('[2010-01-01 14:45, 2010-01-01 15:45)');
ERROR:  conflicting key value violates exclusion constraint "reservation_during_excl"
DETAIL:  Key (during)=(["2010-01-01 14:45:00","2010-01-01 15:45:00"]) conflicts
with existing key (during)=(["2010-01-01 11:30:00","2010-01-01 15:00:00"]).
```

8.18. Domain Types

A *domain* is a user-defined data type that is based on another *underlying type*. Optionally, it can have constraints that restrict its valid values to a subset of what the underlying type would allow. Otherwise it behaves like the underlying type — for example, any operator or function that can be applied to the underlying type will work on the domain type. The underlying type can be any built-in or user-defined base type, enum type, array type, composite type, range type, or another domain.

For example, we could create a domain over integers that accepts only positive integers:

```
CREATE DOMAIN posint AS integer CHECK (VALUE > 0);
CREATE TABLE mytable (id posint);
INSERT INTO mytable VALUES(1);    -- works
INSERT INTO mytable VALUES(-1);  -- fails
```

When an operator or function of the underlying type is applied to a domain value, the domain is automatically down-cast to the underlying type. Thus, for example, the result of `mytable.id - 1` is considered to be of type `integer` not `posint`. We could write `(mytable.id - 1)::posint` to cast the result back to `posint`, causing the domain's constraints to be rechecked. In this case, that would result in an error if the expression had been applied to an `id` value of 1. Assigning a value of the underlying type to a field or variable of the domain type is allowed without writing an explicit cast, but the domain's constraints will be checked.

For additional information see [CREATE DOMAIN](#).

8.19. Object Identifier Types

Object identifiers (OIDs) are used internally by Postgres Pro as primary keys for various system tables. Type `oid` represents an object identifier. There are also several alias types for `oid`, each named *reg-something*. [Table 8.26](#) shows an overview.

The `oid` type is currently implemented as an unsigned four-byte integer. Therefore, it is not large enough to provide database-wide uniqueness in large databases, or even in large individual tables.

The `oid` type itself has few operations beyond comparison. It can be cast to integer, however, and then manipulated using the standard integer operators. (Beware of possible signed-versus-unsigned confusion if you do this.)

The OID alias types have no operations of their own except for specialized input and output routines. These routines are able to accept and display symbolic names for system objects, rather than the raw numeric value that type `oid` would use. The alias types allow simplified lookup of OID values for objects. For example, to examine the `pg_attribute` rows related to a table `mytable`, one could write:

```
SELECT * FROM pg_attribute WHERE attrelid = 'mytable'::regclass;
```

rather than:

```
SELECT * FROM pg_attribute
WHERE attrelid = (SELECT oid FROM pg_class WHERE relname = 'mytable');
```

While that doesn't look all that bad by itself, it's still oversimplified. A far more complicated sub-select would be needed to select the right OID if there are multiple tables named `mytable` in different schemas. The `regclass` input converter handles the table lookup according to the schema path setting, and so it does the “right thing” automatically. Similarly, casting a table's OID to `regclass` is handy for symbolic display of a numeric OID.

Table 8.26. Object Identifier Types

Name	References	Description	Value Example
<code>oid</code>	any	numeric object identifier	564182
<code>regclass</code>	<code>pg_class</code>	relation name	<code>pg_type</code>

Name	References	Description	Value Example
regcollation	pg_collation	collation name	"POSIX"
regconfig	pg_ts_config	text search configura- tion	english
regdictionary	pg_ts_dict	text search dictionary	simple
regnamespace	pg_namespace	namespace name	pg_catalog
regoper	pg_operator	operator name	+
regoperator	pg_operator	operator with argument types	*(integer, integer) or -(NONE, integer)
regproc	pg_proc	function name	sum
regprocedure	pg_proc	function with argument types	sum(int4)
regprofile	pg_profile	profile name	default
regrole	pg_authid	role name	smithee
regtype	pg_type	data type name	integer

All of the OID alias types for objects that are grouped by namespace accept schema-qualified names, and will display schema-qualified names on output if the object would not be found in the current search path without being qualified. For example, `myschema.mytable` is acceptable input for `regclass` (if there is such a table). That value might be output as `myschema.mytable`, or just `mytable`, depending on the current search path. The `regproc` and `regoper` alias types will only accept input names that are unique (not overloaded), so they are of limited use; for most uses `regprocedure` or `regoperator` are more appropriate. For `regoperator`, unary operators are identified by writing `NONE` for the unused operand.

The input functions for these types allow whitespace between tokens, and will fold upper-case letters to lower case, except within double quotes; this is done to make the syntax rules similar to the way object names are written in SQL. Conversely, the output functions will use double quotes if needed to make the output be a valid SQL identifier. For example, the OID of a function named `Foo` (with upper case `F`) taking two integer arguments could be entered as `' "Foo" (int, integer) '::regprocedure`. The output would look like `"Foo" (integer, integer)`. Both the function name and the argument type names could be schema-qualified, too.

Many built-in Postgres Pro functions accept the OID of a table, or another kind of database object, and for convenience are declared as taking `regclass` (or the appropriate OID alias type). This means you do not have to look up the object's OID by hand, but can just enter its name as a string literal. For example, the `nextval(regclass)` function takes a sequence relation's OID, so you could call it like this:

```
nextval('foo')           operates on sequence foo
nextval('FOO')           same as above
nextval('"Foo"')         operates on sequence Foo
nextval('myschema.foo')  operates on myschema.foo
nextval('"myschema".foo') same as above
nextval('foo')           searches search path for foo
```

Note

When you write the argument of such a function as an unadorned literal string, it becomes a constant of type `regclass` (or the appropriate type). Since this is really just an OID, it will track the originally identified object despite later renaming, schema reassignment, etc. This “early binding” behavior is usually desirable for object references in column defaults and views. But sometimes you might want “late binding” where the object reference is resolved at run time. To get late-binding behavior, force the constant to be stored as a `text` constant instead of `regclass`:

```
nextval('foo'::text)      foo is looked up at runtime
```

The `to_regclass()` function and its siblings can also be used to perform run-time lookups. See [Table 9.74](#).

Another practical example of use of `regclass` is to look up the OID of a table listed in the `information_schema` views, which don't supply such OIDs directly. One might for example wish to call the `pg_relation_size()` function, which requires the table OID. Taking the above rules into account, the correct way to do that is

```
SELECT table_schema, table_name,
       pg_relation_size((quote_ident(table_schema) || '.' ||
                           quote_ident(table_name))::regclass)
FROM information_schema.tables
WHERE ...
```

The `quote_ident()` function will take care of double-quoting the identifiers where needed. The seemingly easier

```
SELECT pg_relation_size(table_name)
FROM information_schema.tables
WHERE ...
```

is *not recommended*, because it will fail for tables that are outside your search path or have names that require quoting.

An additional property of most of the OID alias types is the creation of dependencies. If a constant of one of these types appears in a stored expression (such as a column default expression or view), it creates a dependency on the referenced object. For example, if a column has a default expression `nextval('my_seq'::regclass)`, Postgres Pro understands that the default expression depends on the sequence `my_seq`, so the system will not let the sequence be dropped without first removing the default expression. The alternative of `nextval('my_seq'::text)` does not create a dependency. (`regprofile` and `regrole` are exceptions to this property. Constants of these types are not allowed in stored expressions.)

Another identifier type used by the system is `xid`, or transaction (abbreviated `xact`) identifier. This is the data type of the system columns `xmin` and `xmax`. In Postgres Pro Shardman, transaction IDs are implemented as 64-bit counters to prevent transaction ID wraparound. For details, see [Section 24.1.5](#). In some contexts, `xid8` is also used. The `xid` and `xid8` values increase strictly monotonically and cannot be reused in the lifetime of a database cluster. See [Section 67.1](#) for more details.

A third identifier type used by the system is `cid`, or command identifier. This is the data type of the system columns `cmin` and `cmax`. Command identifiers are 32-bit quantities.

A final identifier type used by the system is `tid`, or tuple identifier (row identifier). This is the data type of the system column `ctid`. A tuple ID is a pair (block number, tuple index within block) that identifies the physical location of the row within its table.

(The system columns are further explained in [Section 5.6](#).)

8.20. pg_lsn Type

The `pg_lsn` data type can be used to store LSN (Log Sequence Number) data which is a pointer to a location in the WAL. This type is a representation of `XLogRecPtr` and an internal system type of Postgres Pro.

Internally, an LSN is a 64-bit integer, representing a byte position in the write-ahead log stream. It is printed as two hexadecimal numbers of up to 8 digits each, separated by a slash; for example, `16/B374D848`. The `pg_lsn` type supports the standard comparison operators, like `=` and `>`. Two LSNs can be subtracted using the `-` operator; the result is the number of bytes separating those write-ahead log locations. Also the number of bytes can be added into and subtracted from LSN using the `+` (`pg_lsn, nu-`

`numeric`) and `-(pg_lsn,numeric)` operators, respectively. Note that the calculated LSN should be in the range of `pg_lsn` type, i.e., between 0/0 and FFFFFFFF/FFFFFFF.

8.21. Pseudo-Types

The Postgres Pro type system contains a number of special-purpose entries that are collectively called *pseudo-types*. A pseudo-type cannot be used as a column data type, but it can be used to declare a function's argument or result type. Each of the available pseudo-types is useful in situations where a function's behavior does not correspond to simply taking or returning a value of a specific SQL data type. Table 8.27 lists the existing pseudo-types.

Table 8.27. Pseudo-Types

Name	Description
<code>any</code>	Indicates that a function accepts any input data type.
<code>anyelement</code>	Indicates that a function accepts any data type (see Section 39.2.5).
<code>anyarray</code>	Indicates that a function accepts any array data type (see Section 39.2.5).
<code>anynonarray</code>	Indicates that a function accepts any non-array data type (see Section 39.2.5).
<code>anyenum</code>	Indicates that a function accepts any enum data type (see Section 39.2.5 and Section 8.7).
<code>anyrange</code>	Indicates that a function accepts any range data type (see Section 39.2.5 and Section 8.17).
<code>anymultirange</code>	Indicates that a function accepts any multirange data type (see Section 39.2.5 and Section 8.17).
<code>anycompatible</code>	Indicates that a function accepts any data type, with automatic promotion of multiple arguments to a common data type (see Section 39.2.5).
<code>anycompatiblearray</code>	Indicates that a function accepts any array data type, with automatic promotion of multiple arguments to a common data type (see Section 39.2.5).
<code>anycompatiblenonarray</code>	Indicates that a function accepts any non-array data type, with automatic promotion of multiple arguments to a common data type (see Section 39.2.5).
<code>anycompatiblerange</code>	Indicates that a function accepts any range data type, with automatic promotion of multiple arguments to a common data type (see Section 39.2.5 and Section 8.17).
<code>anycompatiblemultirange</code>	Indicates that a function accepts any multirange data type, with automatic promotion of multiple arguments to a common data type (see Section 39.2.5 and Section 8.17).
<code>cstring</code>	Indicates that a function accepts or returns a null-terminated C string.
<code>internal</code>	Indicates that a function accepts or returns a server-internal data type.
<code>language_handler</code>	A procedural language call handler is declared to return <code>language_handler</code> .
<code>fdw_handler</code>	A foreign-data wrapper handler is declared to return <code>fdw_handler</code> .
<code>table_am_handler</code>	A table access method handler is declared to return <code>table_am_handler</code> .

Name	Description
<code>index_am_handler</code>	An index access method handler is declared to return <code>index_am_handler</code> .
<code>tsm_handler</code>	A tablesample method handler is declared to return <code>tsm_handler</code> .
<code>record</code>	Identifies a function taking or returning an unspecified row type.
<code>trigger</code>	A trigger function is declared to return <code>trigger</code> .
<code>event_trigger</code>	An event trigger function is declared to return <code>event_trigger</code> .
<code>pg_ddl_command</code>	Identifies a representation of DDL commands that is available to event triggers.
<code>void</code>	Indicates that a function returns no value.
<code>unknown</code>	Identifies a not-yet-resolved type, e.g., of an undecorated string literal.

Functions coded in C (whether built-in or dynamically loaded) can be declared to accept or return any of these pseudo-types. It is up to the function author to ensure that the function will behave safely when a pseudo-type is used as an argument type.

Functions coded in procedural languages can use pseudo-types only as allowed by their implementation languages. At present most procedural languages forbid use of a pseudo-type as an argument type, and allow only `void` and `record` as a result type (plus `trigger` or `event_trigger` when the function is used as a trigger or event trigger). Some also support polymorphic functions using the polymorphic pseudo-types, which are shown above and discussed in detail in [Section 39.2.5](#).

The `internal` pseudo-type is used to declare functions that are meant only to be called internally by the database system, and not by direct invocation in an SQL query. If a function has at least one `internal`-type argument then it cannot be called from SQL. To preserve the type safety of this restriction it is important to follow this coding rule: do not create any function that is declared to return `internal` unless it has at least one `internal` argument.

Chapter 9. Functions and Operators

Postgres Pro provides a large number of functions and operators for the built-in data types. This chapter describes most of them, although additional special-purpose functions appear in relevant sections of the manual. Users can also define their own functions and operators, as described in [Part V](#). The `psql` commands `\df` and `\do` can be used to list all available functions and operators, respectively.

The notation used throughout this chapter to describe the argument and result data types of a function or operator is like this:

```
repeat ( text, integer ) → text
```

which says that the function `repeat` takes one text and one integer argument and returns a result of type `text`. The right arrow is also used to indicate the result of an example, thus:

```
repeat ('Pg', 4) → PgPgPgPg
```

If you are concerned about portability then note that most of the functions and operators described in this chapter, with the exception of the most trivial arithmetic and comparison operators and some explicitly marked functions, are not specified by the SQL standard. Some of this extended functionality is present in other SQL database management systems, and in many cases this functionality is compatible and consistent between the various implementations.

9.1. Logical Operators

The usual logical operators are available:

```
boolean AND boolean → boolean
```

```
boolean OR boolean → boolean
```

```
NOT boolean → boolean
```

SQL uses a three-valued logic system with `true`, `false`, and `null`, which represents “unknown”. Observe the following truth tables:

<i>a</i>	<i>b</i>	<i>a AND b</i>	<i>a OR b</i>
TRUE	TRUE	TRUE	TRUE
TRUE	FALSE	FALSE	TRUE
TRUE	NULL	NULL	TRUE
FALSE	FALSE	FALSE	FALSE
FALSE	NULL	FALSE	NULL
NULL	NULL	NULL	NULL

<i>a</i>	NOT <i>a</i>
TRUE	FALSE
FALSE	TRUE
NULL	NULL

The operators `AND` and `OR` are commutative, that is, you can switch the left and right operands without affecting the result. (However, it is not guaranteed that the left operand is evaluated before the right operand. See [Section 4.2.14](#) for more information about the order of evaluation of subexpressions.)

9.2. Comparison Functions and Operators

The usual comparison operators are available, as shown in [Table 9.1](#).

Table 9.1. Comparison Operators

Operator	Description
<code>datatype < datatype → boolean</code>	Less than
<code>datatype > datatype → boolean</code>	Greater than
<code>datatype <= datatype → boolean</code>	Less than or equal to
<code>datatype >= datatype → boolean</code>	Greater than or equal to
<code>datatype = datatype → boolean</code>	Equal
<code>datatype <> datatype → boolean</code>	Not equal
<code>datatype != datatype → boolean</code>	Not equal

Note

`<>` is the standard SQL notation for “not equal”. `!=` is an alias, which is converted to `<>` at a very early stage of parsing. Hence, it is not possible to implement `!=` and `<>` operators that do different things.

These comparison operators are available for all built-in data types that have a natural ordering, including numeric, string, and date/time types. In addition, arrays, composite types, and ranges can be compared if their component data types are comparable.

It is usually possible to compare values of related data types as well; for example `integer > bigint` will work. Some cases of this sort are implemented directly by “cross-type” comparison operators, but if no such operator is available, the parser will coerce the less-general type to the more-general type and apply the latter's comparison operator.

As shown above, all comparison operators are binary operators that return values of type `boolean`. Thus, expressions like `1 < 2 < 3` are not valid (because there is no `<` operator to compare a Boolean value with 3). Use the `BETWEEN` predicates shown below to perform range tests.

There are also some comparison predicates, as shown in [Table 9.2](#). These behave much like operators, but have special syntax mandated by the SQL standard.

Table 9.2. Comparison Predicates

Predicate Description Example(s)
<code>datatype BETWEEN datatype AND datatype → boolean</code> Between (inclusive of the range endpoints). <code>2 BETWEEN 1 AND 3 → t</code> <code>2 BETWEEN 3 AND 1 → f</code>
<code>datatype NOT BETWEEN datatype AND datatype → boolean</code> Not between (the negation of <code>BETWEEN</code>). <code>2 NOT BETWEEN 1 AND 3 → f</code>
<code>datatype BETWEEN SYMMETRIC datatype AND datatype → boolean</code> Between, after sorting the two endpoint values. <code>2 BETWEEN SYMMETRIC 3 AND 1 → t</code>
<code>datatype NOT BETWEEN SYMMETRIC datatype AND datatype → boolean</code> Not between, after sorting the two endpoint values. <code>2 NOT BETWEEN SYMMETRIC 3 AND 1 → f</code>

Predicate	Description	Example(s)
<code>datatype IS DISTINCT FROM datatype</code>	Not equal, treating null as a comparable value.	<code>1 IS DISTINCT FROM NULL → t (rather than NULL)</code> <code>NULL IS DISTINCT FROM NULL → f (rather than NULL)</code>
<code>datatype IS NOT DISTINCT FROM datatype</code>	Equal, treating null as a comparable value.	<code>1 IS NOT DISTINCT FROM NULL → f (rather than NULL)</code> <code>NULL IS NOT DISTINCT FROM NULL → t (rather than NULL)</code>
<code>datatype IS NULL</code>	Test whether value is null.	<code>1.5 IS NULL → f</code>
<code>datatype IS NOT NULL</code>	Test whether value is not null.	<code>'null' IS NOT NULL → t</code>
<code>datatype ISNULL</code>	Test whether value is null (nonstandard syntax).	
<code>datatype NOTNULL</code>	Test whether value is not null (nonstandard syntax).	
<code>boolean IS TRUE</code>	Test whether boolean expression yields true.	<code>true IS TRUE → t</code> <code>NULL::boolean IS TRUE → f (rather than NULL)</code>
<code>boolean IS NOT TRUE</code>	Test whether boolean expression yields false or unknown.	<code>true IS NOT TRUE → f</code> <code>NULL::boolean IS NOT TRUE → t (rather than NULL)</code>
<code>boolean IS FALSE</code>	Test whether boolean expression yields false.	<code>true IS FALSE → f</code> <code>NULL::boolean IS FALSE → f (rather than NULL)</code>
<code>boolean IS NOT FALSE</code>	Test whether boolean expression yields true or unknown.	<code>true IS NOT FALSE → t</code> <code>NULL::boolean IS NOT FALSE → t (rather than NULL)</code>
<code>boolean IS UNKNOWN</code>	Test whether boolean expression yields unknown.	<code>true IS UNKNOWN → f</code> <code>NULL::boolean IS UNKNOWN → t (rather than NULL)</code>
<code>boolean IS NOT UNKNOWN</code>	Test whether boolean expression yields true or false.	<code>true IS NOT UNKNOWN → t</code> <code>NULL::boolean IS NOT UNKNOWN → f (rather than NULL)</code>

The `BETWEEN` predicate simplifies range tests:

```
a BETWEEN x AND y
```

is equivalent to

```
a >= x AND a <= y
```

Notice that `BETWEEN` treats the endpoint values as included in the range. `BETWEEN SYMMETRIC` is like `BETWEEN` except there is no requirement that the argument to the left of `AND` be less than or equal to the argument on the right. If it is not, those two arguments are automatically swapped, so that a nonempty range is always implied.

The various variants of `BETWEEN` are implemented in terms of the ordinary comparison operators, and therefore will work for any data type(s) that can be compared.

Note

The use of `AND` in the `BETWEEN` syntax creates an ambiguity with the use of `AND` as a logical operator. To resolve this, only a limited set of expression types are allowed as the second argument of a `BETWEEN` clause. If you need to write a more complex sub-expression in `BETWEEN`, write parentheses around the sub-expression.

Ordinary comparison operators yield null (signifying “unknown”), not true or false, when either input is null. For example, `7 = NULL` yields null, as does `7 <> NULL`. When this behavior is not suitable, use the `IS [NOT] DISTINCT FROM` predicates:

```
a IS DISTINCT FROM b
```

```
a IS NOT DISTINCT FROM b
```

For non-null inputs, `IS DISTINCT FROM` is the same as the `<>` operator. However, if both inputs are null it returns false, and if only one input is null it returns true. Similarly, `IS NOT DISTINCT FROM` is identical to `=` for non-null inputs, but it returns true when both inputs are null, and false when only one input is null. Thus, these predicates effectively act as though null were a normal data value, rather than “unknown”.

To check whether a value is or is not null, use the predicates:

```
expression IS NULL
```

```
expression IS NOT NULL
```

or the equivalent, but nonstandard, predicates:

```
expression ISNULL
```

```
expression NOTNULL
```

Do *not* write `expression = NULL` because `NULL` is not “equal to” `NULL`. (The null value represents an unknown value, and it is not known whether two unknown values are equal.)

Tip

Some applications might expect that `expression = NULL` returns true if `expression` evaluates to the null value. It is highly recommended that these applications be modified to comply with the SQL standard. However, if that cannot be done the `transform_null_equals` configuration variable is available. If it is enabled, Postgres Pro will convert `x = NULL` clauses to `x IS NULL`.

If the `expression` is row-valued, then `IS NULL` is true when the row expression itself is null or when all the row's fields are null, while `IS NOT NULL` is true when the row expression itself is non-null and all the row's fields are non-null. Because of this behavior, `IS NULL` and `IS NOT NULL` do not always return inverse results for row-valued expressions; in particular, a row-valued expression that contains both null and non-null fields will return false for both tests. For example:

```
SELECT ROW(1,2.5,'this is a test') = ROW(1, 3, 'not the same');
```

```
SELECT ROW(table.*) IS NULL FROM table; -- detect all-null rows
```

```
SELECT ROW(table.*) IS NOT NULL FROM table; -- detect all-non-null rows
```

```
SELECT NOT(ROW(table.*) IS NOT NULL) FROM TABLE; -- detect at least one null in rows
```

In some cases, it may be preferable to write `row IS DISTINCT FROM NULL` or `row IS NOT DISTINCT FROM NULL`, which will simply check whether the overall row value is null without any additional tests on the row fields.

Boolean values can also be tested using the predicates

```
boolean_expression IS TRUE
boolean_expression IS NOT TRUE
boolean_expression IS FALSE
boolean_expression IS NOT FALSE
boolean_expression IS UNKNOWN
boolean_expression IS NOT UNKNOWN
```

These will always return true or false, never a null value, even when the operand is null. A null input is treated as the logical value “unknown”. Notice that `IS UNKNOWN` and `IS NOT UNKNOWN` are effectively the same as `IS NULL` and `IS NOT NULL`, respectively, except that the input expression must be of Boolean type.

Some comparison-related functions are also available, as shown in [Table 9.3](#).

Table 9.3. Comparison Functions

Function	Description	Example(s)
<code>num_nonnulls (VARIADIC "any") → integer</code>	Returns the number of non-null arguments.	<code>num_nonnulls(1, NULL, 2) → 2</code>
<code>num_nulls (VARIADIC "any") → integer</code>	Returns the number of null arguments.	<code>num_nulls(1, NULL, 2) → 1</code>

9.3. Mathematical Functions and Operators

Mathematical operators are provided for many Postgres Pro types. For types without standard mathematical conventions (e.g., date/time types) we describe the actual behavior in subsequent sections.

[Table 9.4](#) shows the mathematical operators that are available for the standard numeric types. Unless otherwise noted, operators shown as accepting *numeric_type* are available for all the types `smallint`, `integer`, `bigint`, `numeric`, `real`, and `double precision`. Operators shown as accepting *integral_type* are available for the types `smallint`, `integer`, and `bigint`. Except where noted, each form of an operator returns the same data type as its argument(s). Calls involving multiple argument data types, such as `integer + numeric`, are resolved by using the type appearing later in these lists.

Table 9.4. Mathematical Operators

Operator	Description	Example(s)
<code>numeric_type + numeric_type → numeric_type</code>	Addition	<code>2 + 3 → 5</code>
<code>+ numeric_type → numeric_type</code>		

Operator	Description	Example(s)
	Unary plus (no operation)	$+ 3.5 \rightarrow 3.5$
	<i>numeric_type</i> - <i>numeric_type</i> → <i>numeric_type</i> Subtraction	$2 - 3 \rightarrow -1$
	- <i>numeric_type</i> → <i>numeric_type</i> Negation	$- (-4) \rightarrow 4$
	<i>numeric_type</i> * <i>numeric_type</i> → <i>numeric_type</i> Multiplication	$2 * 3 \rightarrow 6$
	<i>numeric_type</i> / <i>numeric_type</i> → <i>numeric_type</i> Division (for integral types, division truncates the result towards zero)	$5.0 / 2 \rightarrow 2.5000000000000000$ $5 / 2 \rightarrow 2$ $(-5) / 2 \rightarrow -2$
	<i>numeric_type</i> % <i>numeric_type</i> → <i>numeric_type</i> Modulo (remainder); available for smallint, integer, bigint, and numeric	$5 \% 4 \rightarrow 1$
	<i>numeric</i> ^ <i>numeric</i> → <i>numeric</i> <i>double precision</i> ^ <i>double precision</i> → <i>double precision</i> Exponentiation	$2 ^ 3 \rightarrow 8$ Unlike typical mathematical practice, multiple uses of ^ will associate left to right by default: $2 ^ 3 ^ 3 \rightarrow 512$ $2 ^ (3 ^ 3) \rightarrow 134217728$
	/ <i>double precision</i> → <i>double precision</i> Square root	$ / 25.0 \rightarrow 5$
	/ <i>double precision</i> → <i>double precision</i> Cube root	$ / 64.0 \rightarrow 4$
	@ <i>numeric_type</i> → <i>numeric_type</i> Absolute value	$@ -5.0 \rightarrow 5.0$
	<i>integral_type</i> & <i>integral_type</i> → <i>integral_type</i> Bitwise AND	$91 \& 15 \rightarrow 11$
	<i>integral_type</i> <i>integral_type</i> → <i>integral_type</i> Bitwise OR	$32 3 \rightarrow 35$
	<i>integral_type</i> # <i>integral_type</i> → <i>integral_type</i>	

Operator	Description	Example(s)
	Bitwise exclusive OR	<code>17 # 5 → 20</code>
	<code>~ integral_type → integral_type</code> Bitwise NOT	<code>~1 → -2</code>
	<code>integral_type << integer → integral_type</code> Bitwise shift left	<code>1 << 4 → 16</code>
	<code>integral_type >> integer → integral_type</code> Bitwise shift right	<code>8 >> 2 → 2</code>

Table 9.5 shows the available mathematical functions. Many of these functions are provided in multiple forms with different argument types. Except where noted, any given form of a function returns the same data type as its argument(s); cross-type cases are resolved in the same way as explained above for operators. The functions working with `double precision` data are mostly implemented on top of the host system's C library; accuracy and behavior in boundary cases can therefore vary depending on the host system.

Table 9.5. Mathematical Functions

Function	Description	Example(s)
<code>abs (numeric_type) → numeric_type</code>	Absolute value	<code>abs(-17.4) → 17.4</code>
<code>cbqrt (double precision) → double precision</code>	Cube root	<code>cbqrt(64.0) → 4</code>
<code>ceil (numeric) → numeric</code> <code>ceil (double precision) → double precision</code>	Nearest integer greater than or equal to argument	<code>ceil(42.2) → 43</code> <code>ceil(-42.8) → -42</code>
<code>ceiling (numeric) → numeric</code> <code>ceiling (double precision) → double precision</code>	Nearest integer greater than or equal to argument (same as <code>ceil</code>)	<code>ceiling(95.3) → 96</code>
<code>degrees (double precision) → double precision</code>	Converts radians to degrees	<code>degrees(0.5) → 28.64788975654116</code>
<code>div (y numeric, x numeric) → numeric</code>	Integer quotient of y/x (truncates towards zero)	<code>div(9, 4) → 2</code>
<code>erf (double precision) → double precision</code>		

Function	Description	Example(s)
Error function		<code>erf(1.0) → 0.8427007929497149</code>
erfc (double precision) → double precision	Complementary error function (1 - erf(x) , without loss of precision for large inputs)	<code>erfc(1.0) → 0.15729920705028513</code>
exp (numeric) → numeric exp (double precision) → double precision	Exponential (e raised to the given power)	<code>exp(1.0) → 2.7182818284590452</code>
factorial (bigint) → numeric	Factorial	<code>factorial(5) → 120</code>
floor (numeric) → numeric floor (double precision) → double precision	Nearest integer less than or equal to argument	<code>floor(42.8) → 42</code> <code>floor(-42.8) → -43</code>
gcd (numeric_type , numeric_type) → numeric_type	Greatest common divisor (the largest positive number that divides both inputs with no remainder); returns 0 if both inputs are zero; available for integer, bigint, and numeric	<code>gcd(1071, 462) → 21</code>
lcm (numeric_type , numeric_type) → numeric_type	Least common multiple (the smallest strictly positive number that is an integral multiple of both inputs); returns 0 if either input is zero; available for integer, bigint, and numeric	<code>lcm(1071, 462) → 23562</code>
ln (numeric) → numeric ln (double precision) → double precision	Natural logarithm	<code>ln(2.0) → 0.6931471805599453</code>
log (numeric) → numeric log (double precision) → double precision	Base 10 logarithm	<code>log(100) → 2</code>
log10 (numeric) → numeric log10 (double precision) → double precision	Base 10 logarithm (same as log)	<code>log10(1000) → 3</code>
log (b numeric, x numeric) → numeric	Logarithm of x to base b	<code>log(2.0, 64.0) → 6.0000000000000000</code>
min_scale (numeric) → integer	Minimum scale (number of fractional decimal digits) needed to represent the supplied value precisely	

Function	Description	Example(s)
<code>min_scale(8.4100)</code>		<code>→ 2</code>
<code>mod (y numeric_type , x numeric_type)</code>	<code>→ numeric_type</code> Remainder of y/x ; available for <code>smallint</code> , <code>integer</code> , <code>bigint</code> , and <code>numeric</code>	<code>mod(9, 4) → 1</code>
<code>pi ()</code>	<code>→ double precision</code> Approximate value of π	<code>pi () → 3.141592653589793</code>
<code>power (a numeric , b numeric)</code>	<code>→ numeric</code>	
<code>power (a double precision , b double precision)</code>	<code>→ double precision</code> a raised to the power of b	<code>power(9, 3) → 729</code>
<code>radians (double precision)</code>	<code>→ double precision</code> Converts degrees to radians	<code>radians(45.0) → 0.7853981633974483</code>
<code>round (numeric)</code>	<code>→ numeric</code>	
<code>round (double precision)</code>	<code>→ double precision</code> Rounds to nearest integer. For <code>numeric</code> , ties are broken by rounding away from zero. For <code>double precision</code> , the tie-breaking behavior is platform dependent, but “round to nearest even” is the most common rule.	<code>round(42.4) → 42</code>
<code>round (v numeric , s integer)</code>	<code>→ numeric</code> Rounds v to s decimal places. Ties are broken by rounding away from zero.	<code>round(42.4382, 2) → 42.44</code> <code>round(1234.56, -1) → 1230</code>
<code>scale (numeric)</code>	<code>→ integer</code> Scale of the argument (the number of decimal digits in the fractional part)	<code>scale(8.4100) → 4</code>
<code>sign (numeric)</code>	<code>→ numeric</code>	
<code>sign (double precision)</code>	<code>→ double precision</code> Sign of the argument (-1, 0, or +1)	<code>sign(-8.4) → -1</code>
<code>sqrt (numeric)</code>	<code>→ numeric</code>	
<code>sqrt (double precision)</code>	<code>→ double precision</code> Square root	<code>sqrt(2) → 1.4142135623730951</code>
<code>trim_scale (numeric)</code>	<code>→ numeric</code> Reduces the value's scale (number of fractional decimal digits) by removing trailing zeroes	<code>trim_scale(8.4100) → 8.41</code>
<code>trunc (numeric)</code>	<code>→ numeric</code>	
<code>trunc (double precision)</code>	<code>→ double precision</code> Truncates to integer (towards zero)	<code>trunc(42.8) → 42</code>

Function	Description	Example(s)
		<code>trunc(-42.8) → -42</code>
	<code>trunc (v numeric, s integer) → numeric</code> Truncates <i>v</i> to <i>s</i> decimal places	<code>trunc(42.4382, 2) → 42.43</code>
	<code>width_bucket (operand numeric, low numeric, high numeric, count integer) → integer</code> <code>width_bucket (operand double precision, low double precision, high double precision, count integer) → integer</code> Returns the number of the bucket in which <i>operand</i> falls in a histogram having <i>count</i> equal-width buckets spanning the range <i>low</i> to <i>high</i> . Returns 0 or <i>count</i> +1 for an input outside that range.	<code>width_bucket(5.35, 0.024, 10.06, 5) → 3</code>
	<code>width_bucket (operand anycompatible, thresholds anycompatiblearray) → integer</code> Returns the number of the bucket in which <i>operand</i> falls given an array listing the lower bounds of the buckets. Returns 0 for an input less than the first lower bound. <i>operand</i> and the array elements can be of any type having standard comparison operators. The <i>thresholds</i> array <i>must be sorted</i> , smallest first, or unexpected results will be obtained.	<code>width_bucket(now(), array['yesterday', 'today', 'tomorrow']::time-stampztz[]) → 2</code>

Table 9.6 shows functions for generating random numbers.

Table 9.6. Random Functions

Function	Description	Example(s)
	<code>random () → double precision</code> Returns a random value in the range $0.0 \leq x < 1.0$	<code>random() → 0.897124072839091</code>
	<code>random (min integer, max integer) → integer</code> <code>random (min bigint, max bigint) → bigint</code> <code>random (min numeric, max numeric) → numeric</code> Returns a random value in the range $min \leq x \leq max$. For type <i>numeric</i> , the result will have the same number of fractional decimal digits as <i>min</i> or <i>max</i> , whichever has more.	<code>random(1, 10) → 7</code> <code>random(-0.499, 0.499) → 0.347</code>
	<code>random_normal ([mean double precision [, stddev double precision]]) → double precision</code> Returns a random value from the normal distribution with the given parameters; <i>mean</i> defaults to 0.0 and <i>stddev</i> defaults to 1.0	<code>random_normal(0.0, 1.0) → 0.051285419</code>
	<code>setseed (double precision) → void</code> Sets the seed for subsequent <code>random()</code> and <code>random_normal()</code> calls; argument must be between -1.0 and 1.0, inclusive	<code>setseed(0.12345)</code>

The `random()` and `random_normal()` functions listed in Table 9.6 use a deterministic pseudo-random number generator. It is fast but not suitable for cryptographic applications; see the [pgcrypto](#) module for a

more secure alternative. If `setseed()` is called, the series of results of subsequent calls to these functions in the current session can be repeated by re-issuing `setseed()` with the same argument. Without any prior `setseed()` call in the same session, the first call to any of these functions obtains a seed from a platform-dependent source of random bits.

[Table 9.7](#) shows the available trigonometric functions. Each of these functions comes in two variants, one that measures angles in radians and one that measures angles in degrees.

Table 9.7. Trigonometric Functions

Function Description Example(s)
<code>acos (double precision) → double precision</code> Inverse cosine, result in radians <code>acos(1) → 0</code>
<code>acosd (double precision) → double precision</code> Inverse cosine, result in degrees <code>acosd(0.5) → 60</code>
<code>asin (double precision) → double precision</code> Inverse sine, result in radians <code>asin(1) → 1.5707963267948966</code>
<code>asind (double precision) → double precision</code> Inverse sine, result in degrees <code>asind(0.5) → 30</code>
<code>atan (double precision) → double precision</code> Inverse tangent, result in radians <code>atan(1) → 0.7853981633974483</code>
<code>atand (double precision) → double precision</code> Inverse tangent, result in degrees <code>atand(1) → 45</code>
<code>atan2 (y double precision, x double precision) → double precision</code> Inverse tangent of y/x , result in radians <code>atan2(1, 0) → 1.5707963267948966</code>
<code>atan2d (y double precision, x double precision) → double precision</code> Inverse tangent of y/x , result in degrees <code>atan2d(1, 0) → 90</code>
<code>cos (double precision) → double precision</code> Cosine, argument in radians <code>cos(0) → 1</code>
<code>cosd (double precision) → double precision</code> Cosine, argument in degrees <code>cosd(60) → 0.5</code>
<code>cot (double precision) → double precision</code> Cotangent, argument in radians <code>cot(0.5) → 1.830487721712452</code>
<code>cotd (double precision) → double precision</code> Cotangent, argument in degrees

Function	Description	Example(s)
		<code>cotd(45) → 1</code>
<code>sin (double precision)</code>	→ double precision Sine, argument in radians	<code>sin(1) → 0.8414709848078965</code>
<code>sind (double precision)</code>	→ double precision Sine, argument in degrees	<code>sind(30) → 0.5</code>
<code>tan (double precision)</code>	→ double precision Tangent, argument in radians	<code>tan(1) → 1.5574077246549023</code>
<code>tand (double precision)</code>	→ double precision Tangent, argument in degrees	<code>tand(45) → 1</code>

Note

Another way to work with angles measured in degrees is to use the unit transformation functions `radians()` and `degrees()` shown earlier. However, using the degree-based trigonometric functions is preferred, as that way avoids round-off error for special cases such as `sind(30)`.

Table 9.8 shows the available hyperbolic functions.

Table 9.8. Hyperbolic Functions

Function	Description	Example(s)
<code>sinh (double precision)</code>	→ double precision Hyperbolic sine	<code>sinh(1) → 1.1752011936438014</code>
<code>cosh (double precision)</code>	→ double precision Hyperbolic cosine	<code>cosh(0) → 1</code>
<code>tanh (double precision)</code>	→ double precision Hyperbolic tangent	<code>tanh(1) → 0.7615941559557649</code>
<code>asinh (double precision)</code>	→ double precision Inverse hyperbolic sine	<code>asinh(1) → 0.881373587019543</code>
<code>acosh (double precision)</code>	→ double precision Inverse hyperbolic cosine	<code>acosh(1) → 0</code>
<code>atanh (double precision)</code>	→ double precision Inverse hyperbolic tangent	<code>atanh(0.5) → 0.5493061443340548</code>

9.4. String Functions and Operators

This section describes functions and operators for examining and manipulating string values. Strings in this context include values of the types `character`, `character varying`, and `text`. Except where noted, these functions and operators are declared to accept and return type `text`. They will interchangeably accept `character varying` arguments. Values of type `character` will be converted to `text` before the function or operator is applied, resulting in stripping any trailing spaces in the `character` value.

SQL defines some string functions that use key words, rather than commas, to separate arguments. Details are in [Table 9.9](#). Postgres Pro also provides versions of these functions that use the regular function invocation syntax (see [Table 9.10](#)).

Note

The string concatenation operator (`||`) will accept non-string input, so long as at least one input is of string type, as shown in [Table 9.9](#). For other cases, inserting an explicit coercion to `text` can be used to have non-string input accepted.

Table 9.9. SQL String Functions and Operators

Function/Operator	Description	Example(s)
<code>text text</code>	<code>→ text</code> Concatenates the two strings.	<code>'Post' 'greSQL' → PostgreSQL</code>
<code>text anynonarray</code> <code>anynonarray text</code>	<code>→ text</code> Converts the non-string input to text, then concatenates the two strings. (The non-string input cannot be of an array type, because that would create ambiguity with the <code>array </code> operators. If you want to concatenate an array's text equivalent, cast it to <code>text</code> explicitly.)	<code>'Value: ' 42 → Value: 42</code>
<code>btrim (string text [, characters text])</code>	<code>→ text</code> Removes the longest string containing only characters in <i>characters</i> (a space by default) from the start and end of <i>string</i> .	<code>btrim('xyxtrimyyx', 'xyz') → trim</code>
<code>text IS [NOT] [form] NORMALIZED</code>	<code>→ boolean</code> Checks whether the string is in the specified Unicode normalization form. The optional <i>form</i> key word specifies the form: <code>NFC</code> (the default), <code>NFD</code> , <code>NFKC</code> , or <code>NFKD</code> . This expression can only be used when the server encoding is <code>UTF8</code> . Note that checking for normalization using this expression is often faster than normalizing possibly already normalized strings.	<code>U&'\0061\0308bc' IS NFD NORMALIZED → t</code>
<code>bit_length (text)</code>	<code>→ integer</code> Returns number of bits in the string (8 times the <code>octet_length</code>).	<code>bit_length('jose') → 32</code>
<code>char_length (text)</code> <code>character_length (text)</code>	<code>→ integer</code> Returns number of characters in the string.	<code>char_length('josé') → 4</code>
<code>lower (text)</code>	<code>→ text</code> Converts the string to all lower case, according to the rules of the database's locale.	

Function/Operator	Description	Example(s)
		<code>lower('TOM') → tom</code>
<code>lpad (string text, length integer [, fill text]) → text</code>	Extends the <i>string</i> to length <i>length</i> by prepending the characters <i>fill</i> (a space by default). If the <i>string</i> is already longer than <i>length</i> then it is truncated (on the right).	<code>lpad('hi', 5, 'xy') → xyxhi</code>
<code>ltrim (string text [, characters text]) → text</code>	Removes the longest string containing only characters in <i>characters</i> (a space by default) from the start of <i>string</i> .	<code>ltrim('zzzytest', 'xyz') → test</code>
<code>normalize (text [, form]) → text</code>	Converts the string to the specified Unicode normalization form. The optional <i>form</i> key word specifies the form: NFC (the default), NFD, NFKC, or NFKD. This function can only be used when the server encoding is UTF8.	<code>normalize(U&'\0061\0308bc', NFC) → U&'\00E4bc'</code>
<code>octet_length (text) → integer</code>	Returns number of bytes in the string.	<code>octet_length('josé') → 5 (if server encoding is UTF8)</code>
<code>octet_length (character) → integer</code>	Returns number of bytes in the string. Since this version of the function accepts type <i>character</i> directly, it will not strip trailing spaces.	<code>octet_length('abc '::character(4)) → 4</code>
<code>overlay (string text PLACING newsubstring text FROM start integer [FOR count integer]) → text</code>	Replaces the substring of <i>string</i> that starts at the <i>start</i> 'th character and extends for <i>count</i> characters with <i>newsubstring</i> . If <i>count</i> is omitted, it defaults to the length of <i>newsubstring</i> .	<code>overlay('Txxxxas' placing 'hom' from 2 for 4) → Thomas</code>
<code>position (substring text IN string text) → integer</code>	Returns first starting index of the specified <i>substring</i> within <i>string</i> , or zero if it's not present.	<code>position('om' in 'Thomas') → 3</code>
<code>rpadd (string text, length integer [, fill text]) → text</code>	Extends the <i>string</i> to length <i>length</i> by appending the characters <i>fill</i> (a space by default). If the <i>string</i> is already longer than <i>length</i> then it is truncated.	<code>rpadd('hi', 5, 'xy') → hixyx</code>
<code>rtrim (string text [, characters text]) → text</code>	Removes the longest string containing only characters in <i>characters</i> (a space by default) from the end of <i>string</i> .	<code>rtrim('testxxxz', 'xyz') → test</code>
<code>substring (string text [FROM start integer] [FOR count integer]) → text</code>	Extracts the substring of <i>string</i> starting at the <i>start</i> 'th character if that is specified, and stopping after <i>count</i> characters if that is specified. Provide at least one of <i>start</i> and <i>count</i> .	<code>substring('Thomas' from 2 for 3) → hom</code> <code>substring('Thomas' from 3) → omas</code> <code>substring('Thomas' for 2) → Th</code>

Function/Operator	Description	Example(s)
<code>substring (<i>string</i> text FROM <i>pattern</i> text)</code>	Extracts the first substring matching POSIX regular expression; see Section 9.7.3 .	<code>substring('Thomas' from '...\$') → mas</code>
<code>substring (<i>string</i> text SIMILAR <i>pattern</i> text ESCAPE <i>escape</i> text)</code> <code>substring (<i>string</i> text FROM <i>pattern</i> text FOR <i>escape</i> text)</code>	Extracts the first substring matching SQL regular expression; see Section 9.7.2 . The first form has been specified since SQL:2003; the second form was only in SQL:1999 and should be considered obsolete.	<code>substring('Thomas' similar '%#"o_a#"' escape '#') → oma</code>
<code>trim ([LEADING TRAILING BOTH] [<i>characters</i> text] FROM <i>string</i> text)</code>	Removes the longest string containing only characters in <i>characters</i> (a space by default) from the start, end, or both ends (BOTH is the default) of <i>string</i> .	<code>trim(both 'xyz' from 'yxTomxx') → Tom</code>
<code>trim ([LEADING TRAILING BOTH] [FROM] <i>string</i> text [, <i>characters</i> text])</code>	This is a non-standard syntax for <code>trim()</code> .	<code>trim(both from 'yxTomxx', 'xyz') → Tom</code>
<code>unicode_assigned (text)</code>	Returns <code>true</code> if all characters in the string are assigned Unicode codepoints; <code>false</code> otherwise. This function can only be used when the server encoding is UTF8.	
<code>upper (text)</code>	Converts the string to all upper case, according to the rules of the database's locale.	<code>upper('tom') → TOM</code>

Additional string manipulation functions and operators are available and are listed in [Table 9.10](#). (Some of these are used internally to implement the SQL-standard string functions listed in [Table 9.9](#).) There are also pattern-matching operators, which are described in [Section 9.7](#), and operators for full-text search, which are described in [Chapter 12](#).

Table 9.10. Other String Functions and Operators

Function/Operator	Description	Example(s)
<code>text ^@ text</code>	Returns <code>true</code> if the first string starts with the second string (equivalent to the <code>starts_with()</code> function).	<code>'alphabet' ^@ 'alph' → t</code>
<code>ascii (text)</code>	Returns the numeric code of the first character of the argument. In UTF8 encoding, returns the Unicode code point of the character. In other multibyte encodings, the argument must be an ASCII character.	<code>ascii('x') → 120</code>
<code>chr (integer)</code>	Returns the character with the given code. In UTF8 encoding the argument is treated as a Unicode code point. In other multibyte encodings the argument must designate an ASCII character. <code>chr(0)</code> is disallowed because text data types cannot store that character.	<code>chr(65) → A</code>

Function/Operator	Description	Example(s)
<code>concat</code>	<code>concat (val1 "any" [, val2 "any" [, ...]]) → text</code> Concatenates the text representations of all the arguments. NULL arguments are ignored.	<code>concat('abcde', 2, NULL, 22) → abcde222</code>
<code>concat_ws</code>	<code>concat_ws (sep text, val1 "any" [, val2 "any" [, ...]]) → text</code> Concatenates all but the first argument, with separators. The first argument is used as the separator string, and should not be NULL. Other NULL arguments are ignored.	<code>concat_ws(',', 'abcde', 2, NULL, 22) → abcde,2,22</code>
<code>format</code>	<code>format (formatstr text [, formatarg "any" [, ...]]) → text</code> Formats arguments according to a format string; see Section 9.4.1 . This function is similar to the C function <code>sprintf</code> .	<code>format('Hello %s, %1\$s', 'World') → Hello World, World</code>
<code>initcap</code>	<code>initcap (text) → text</code> Converts the first letter of each word to upper case and the rest to lower case. Words are sequences of alphanumeric characters separated by non-alphanumeric characters.	<code>initcap('hi THOMAS') → Hi Thomas</code>
<code>left</code>	<code>left (string text, n integer) → text</code> Returns first <i>n</i> characters in the string, or when <i>n</i> is negative, returns all but last <i> n </i> characters.	<code>left('abcde', 2) → ab</code>
<code>length</code>	<code>length (text) → integer</code> Returns the number of characters in the string.	<code>length('jose') → 4</code>
<code>md5</code>	<code>md5 (text) → text</code> Computes the MD5 hash of the argument, with the result written in hexadecimal.	<code>md5('abc') → 900150983cd24fb0d6963f7d28e17f72</code>
<code>parse_ident</code>	<code>parse_ident (qualified_identifier text [, strict_mode boolean DEFAULT true]) → text[]</code> Splits <i>qualified_identifier</i> into an array of identifiers, removing any quoting of individual identifiers. By default, extra characters after the last identifier are considered an error; but if the second parameter is <code>false</code> , then such extra characters are ignored. (This behavior is useful for parsing names for objects like functions.) Note that this function does not truncate over-length identifiers. If you want truncation you can cast the result to <code>name[]</code> .	<code>parse_ident('"SomeSchema".someTable') → {SomeSchema,sometable}</code>
<code>pg_client_encoding</code>	<code>pg_client_encoding () → name</code> Returns current client encoding name.	<code>pg_client_encoding() → UTF8</code>
<code>quote_ident</code>	<code>quote_ident (text) → text</code> Returns the given string suitably quoted to be used as an identifier in an SQL statement string. Quotes are added only if necessary (i.e., if the string contains non-identifier characters or would be case-folded). Embedded quotes are properly doubled. See also Example 44.1 .	<code>quote_ident('Foo bar') → "Foo bar"</code>
<code>quote_literal</code>	<code>quote_literal (text) → text</code> Returns the given string suitably quoted to be used as a string literal in an SQL statement string. Embedded single-quotes and backslashes are properly doubled. Note that <code>quote_literal</code> returns null on null input; if the argument might be null, <code>quote_nullable</code> is often more suitable. See also Example 44.1 .	

Function/Operator	Description	Example(s)
		<code>quote_literal(E'O\Reilly')</code> → <code>'O'Reilly'</code>
<code>quote_literal</code>	<code>(anyelement) → text</code> Converts the given value to text and then quotes it as a literal. Embedded single-quotes and backslashes are properly doubled.	<code>quote_literal(42.5)</code> → <code>'42.5'</code>
<code>quote_nullable</code>	<code>(text) → text</code> Returns the given string suitably quoted to be used as a string literal in an SQL statement string; or, if the argument is null, returns <code>NULL</code> . Embedded single-quotes and backslashes are properly doubled. See also Example 44.1 .	<code>quote_nullable(NULL)</code> → <code>NULL</code>
<code>quote_nullable</code>	<code>(anyelement) → text</code> Converts the given value to text and then quotes it as a literal; or, if the argument is null, returns <code>NULL</code> . Embedded single-quotes and backslashes are properly doubled.	<code>quote_nullable(42.5)</code> → <code>'42.5'</code>
<code>regexp_count</code>	<code>(string text, pattern text [, start integer [, flags text]]) → integer</code> Returns the number of times the POSIX regular expression <i>pattern</i> matches in the <i>string</i> ; see Section 9.7.3 .	<code>regexp_count('123456789012', '\d\d\d', 2)</code> → 3
<code>regexp_instr</code>	<code>(string text, pattern text [, start integer [, N integer [, endoption integer [, flags text [, subexpr integer]]]]]) → integer</code> Returns the position within <i>string</i> where the <i>N</i> 'th match of the POSIX regular expression <i>pattern</i> occurs, or zero if there is no such match; see Section 9.7.3 .	<code>regexp_instr('ABCDEF', 'c(..)', 1, 1, 0, 'i')</code> → 3 <code>regexp_instr('ABCDEF', 'c(..)', 1, 1, 0, 'i', 2)</code> → 5
<code>regexp_like</code>	<code>(string text, pattern text [, flags text]) → boolean</code> Checks whether a match of the POSIX regular expression <i>pattern</i> occurs within <i>string</i> ; see Section 9.7.3 .	<code>regexp_like('Hello World', 'world\$', 'i')</code> → <code>t</code>
<code>regexp_match</code>	<code>(string text, pattern text [, flags text]) → text[]</code> Returns substrings within the first match of the POSIX regular expression <i>pattern</i> to the <i>string</i> ; see Section 9.7.3 .	<code>regexp_match('foobarbequebaz', '(bar) (beque)')</code> → {bar, beque}
<code>regexp_matches</code>	<code>(string text, pattern text [, flags text]) → setof text[]</code> Returns substrings within the first match of the POSIX regular expression <i>pattern</i> to the <i>string</i> , or substrings within all such matches if the <i>g</i> flag is used; see Section 9.7.3 .	<code>regexp_matches('foobarbequebaz', 'ba.', 'g')</code> → {bar} {baz}
<code>regexp_replace</code>	<code>(string text, pattern text, replacement text [, start integer] [, flags text]) → text</code> Replaces the substring that is the first match to the POSIX regular expression <i>pattern</i> , or all such matches if the <i>g</i> flag is used; see Section 9.7.3 .	<code>regexp_replace('Thomas', '[mN]a.', 'M')</code> → ThM

Function/Operator	Description	Example(s)
<code>regexp_replace</code>	(<i>string</i> text, <i>pattern</i> text, <i>replacement</i> text, <i>start</i> integer, <i>N</i> integer [, <i>flags</i> text]) → text Replaces the substring that is the <i>N</i> 'th match to the POSIX regular expression <i>pattern</i> , or all such matches if <i>N</i> is zero; see Section 9.7.3 .	<code>regexp_replace('Thomas', '.', 'X', 3, 2)</code> → ThoXas
<code>regexp_split_to_array</code>	(<i>string</i> text, <i>pattern</i> text [, <i>flags</i> text]) → text [] Splits <i>string</i> using a POSIX regular expression as the delimiter, producing an array of results; see Section 9.7.3 .	<code>regexp_split_to_array('hello world', '\s+')</code> → {hello, world}
<code>regexp_split_to_table</code>	(<i>string</i> text, <i>pattern</i> text [, <i>flags</i> text]) → set of text Splits <i>string</i> using a POSIX regular expression as the delimiter, producing a set of results; see Section 9.7.3 .	<code>regexp_split_to_table('hello world', '\s+')</code> → hello world
<code>regexp_substr</code>	(<i>string</i> text, <i>pattern</i> text [, <i>start</i> integer [, <i>N</i> integer [, <i>flags</i> text [, <i>subexpr</i> integer]]]]) → text Returns the substring within <i>string</i> that matches the <i>N</i> 'th occurrence of the POSIX regular expression <i>pattern</i> , or NULL if there is no such match; see Section 9.7.3 .	<code>regexp_substr('ABCDEF', 'c(.)(..)', 1, 1, 'i')</code> → CDEF <code>regexp_substr('ABCDEF', 'c(.)(..)', 1, 1, 'i', 2)</code> → EF
<code>repeat</code>	(<i>string</i> text, <i>number</i> integer) → text Repeats <i>string</i> the specified <i>number</i> of times.	<code>repeat('Pg', 4)</code> → PgPgPgPg
<code>replace</code>	(<i>string</i> text, <i>from</i> text, <i>to</i> text) → text Replaces all occurrences in <i>string</i> of substring <i>from</i> with substring <i>to</i> .	<code>replace('abcdefabcdef', 'cd', 'XX')</code> → abXXefabXXef
<code>reverse</code>	(text) → text Reverses the order of the characters in the string.	<code>reverse('abcde')</code> → edcba
<code>right</code>	(<i>string</i> text, <i>n</i> integer) → text Returns last <i>n</i> characters in the string, or when <i>n</i> is negative, returns all but first <i>n</i> characters.	<code>right('abcde', 2)</code> → de
<code>split_part</code>	(<i>string</i> text, <i>delimiter</i> text, <i>n</i> integer) → text Splits <i>string</i> at occurrences of <i>delimiter</i> and returns the <i>n</i> 'th field (counting from one), or when <i>n</i> is negative, returns the <i>n</i> "th-from-last field.	<code>split_part('abc~@~def~@~ghi', '~@~', 2)</code> → def <code>split_part('abc,def,ghi,jkl', ',', -2)</code> → ghi
<code>starts_with</code>	(<i>string</i> text, <i>prefix</i> text) → boolean Returns true if <i>string</i> starts with <i>prefix</i> .	<code>starts_with('alphabet', 'alph')</code> → t
<code>string_to_array</code>	(<i>string</i> text, <i>delimiter</i> text [, <i>null_string</i> text]) → text []	

Function/Operator	Description	Example(s)
	Splits the <i>string</i> at occurrences of <i>delimiter</i> and forms the resulting fields into a text array. If <i>delimiter</i> is NULL, each character in the <i>string</i> will become a separate element in the array. If <i>delimiter</i> is an empty string, then the <i>string</i> is treated as a single field. If <i>null_string</i> is supplied and is not NULL, fields matching that string are replaced by NULL. See also array_to_string .	<code>string_to_array('xx~yy~zz', '~', 'yy') → {xx, NULL, zz}</code>
<code>string_to_table (string text, delimiter text [, null_string text]) → set of text</code>	Splits the <i>string</i> at occurrences of <i>delimiter</i> and returns the resulting fields as a set of text rows. If <i>delimiter</i> is NULL, each character in the <i>string</i> will become a separate row of the result. If <i>delimiter</i> is an empty string, then the <i>string</i> is treated as a single field. If <i>null_string</i> is supplied and is not NULL, fields matching that string are replaced by NULL.	<code>string_to_table('xx~^yy~^zz', '~^', 'yy') →</code> <code>xx</code> <code>NULL</code> <code>zz</code>
<code>strpos (string text, substring text) → integer</code>	Returns first starting index of the specified <i>substring</i> within <i>string</i> , or zero if it's not present. (Same as <code>position(substring in string)</code> , but note the reversed argument order.)	<code>strpos('high', 'ig') → 2</code>
<code>substr (string text, start integer [, count integer]) → text</code>	Extracts the substring of <i>string</i> starting at the <i>start</i> 'th character, and extending for <i>count</i> characters if that is specified. (Same as <code>substring(string from start for count)</code> .)	<code>substr('alphabet', 3) → phabet</code> <code>substr('alphabet', 3, 2) → ph</code>
<code>to_ascii (string text) → text</code> <code>to_ascii (string text, encoding name) → text</code> <code>to_ascii (string text, encoding integer) → text</code>	Converts <i>string</i> to ASCII from another encoding, which may be identified by name or number. If <i>encoding</i> is omitted the database encoding is assumed (which in practice is the only useful case). The conversion consists primarily of dropping accents. Conversion is only supported from LATIN1, LATIN2, LATIN9, and WIN1250 encodings.	<code>to_ascii('Karél') → Karel</code>
<code>to_bin (integer) → text</code> <code>to_bin (bigint) → text</code>	Converts the number to its equivalent two's complement binary representation.	<code>to_bin(2147483647) → 11111111111111111111111111111111</code> <code>to_bin(-1234) → 1111111111111111111111111111111101100101110</code>
<code>to_hex (integer) → text</code> <code>to_hex (bigint) → text</code>	Converts the number to its equivalent two's complement hexadecimal representation.	<code>to_hex(2147483647) → 7fffffff</code> <code>to_hex(-1234) → fffffb2e</code>
<code>to_oct (integer) → text</code> <code>to_oct (bigint) → text</code>		

Function/Operator	Description	Example(s)
	Converts the number to its equivalent two's complement octal representation.	<code>to_oct (2147483647) → 17777777777</code> <code>to_oct (-1234) → 37777775456</code>
<code>translate (string text, from text, to text)</code>	Replaces each character in <i>string</i> that matches a character in the <i>from</i> set with the corresponding character in the <i>to</i> set. If <i>from</i> is longer than <i>to</i> , occurrences of the extra characters in <i>from</i> are deleted.	<code>translate('12345', '143', 'ax') → a2x5</code>
<code>unistr (text)</code>	Evaluate escaped Unicode characters in the argument. Unicode characters can be specified as <code>\XXXX</code> (4 hexadecimal digits), <code>\+XXXXXX</code> (6 hexadecimal digits), <code>\uXXXX</code> (4 hexadecimal digits), or <code>\UXXXXXXXX</code> (8 hexadecimal digits). To specify a backslash, write two backslashes. All other characters are taken literally. If the server encoding is not UTF-8, the Unicode code point identified by one of these escape sequences is converted to the actual server encoding; an error is reported if that's not possible. This function provides a (non-standard) alternative to string constants with Unicode escapes (see Section 4.1.2.3).	<code>unistr('d\0061t\+000061') → data</code> <code>unistr('d\u0061t\U00000061') → data</code>

The `concat`, `concat_ws` and `format` functions are variadic, so it is possible to pass the values to be concatenated or formatted as an array marked with the `VARIADIC` keyword (see [Section 39.5.6](#)). The array's elements are treated as if they were separate ordinary arguments to the function. If the variadic array argument is `NULL`, `concat` and `concat_ws` return `NULL`, but `format` treats a `NULL` as a zero-element array.

See also the aggregate function `string_agg` in [Section 9.21](#), and the functions for converting between strings and the `bytea` type in [Table 9.13](#).

9.4.1. format

The function `format` produces output formatted according to a format string, in a style similar to the C function `printf`.

```
format(formatstr text [, formatarg "any" [, ...] ])
```

formatstr is a format string that specifies how the result should be formatted. Text in the format string is copied directly to the result, except where *format specifiers* are used. Format specifiers act as placeholders in the string, defining how subsequent function arguments should be formatted and inserted into the result. Each *formatarg* argument is converted to text according to the usual output rules for its data type, and then formatted and inserted into the result string according to the format specifier(s).

Format specifiers are introduced by a `%` character and have the form

```
%[position] [flags] [width] type
```

where the component fields are:

position (optional)

A string of the form `n$` where *n* is the index of the argument to print. Index 1 means the first argument after *formatstr*. If the *position* is omitted, the default is to use the next argument in sequence.

flags (optional)

Additional options controlling how the format specifier's output is formatted. Currently the only supported flag is a minus sign (-) which will cause the format specifier's output to be left-justified. This has no effect unless the *width* field is also specified.

width (optional)

Specifies the *minimum* number of characters to use to display the format specifier's output. The output is padded on the left or right (depending on the - flag) with spaces as needed to fill the width. A too-small width does not cause truncation of the output, but is simply ignored. The width may be specified using any of the following: a positive integer; an asterisk (*) to use the next function argument as the width; or a string of the form *n\$ to use the *n*th function argument as the width.

If the width comes from a function argument, that argument is consumed before the argument that is used for the format specifier's value. If the width argument is negative, the result is left aligned (as if the - flag had been specified) within a field of length *abs(width)*.

type (required)

The type of format conversion to use to produce the format specifier's output. The following types are supported:

- **s** formats the argument value as a simple string. A null value is treated as an empty string.
- **I** treats the argument value as an SQL identifier, double-quoting it if necessary. It is an error for the value to be null (equivalent to `quote_ident`).
- **L** quotes the argument value as an SQL literal. A null value is displayed as the string `NULL`, without quotes (equivalent to `quote_nullable`).

In addition to the format specifiers described above, the special sequence %% may be used to output a literal % character.

Here are some examples of the basic format conversions:

```
SELECT format('Hello %s', 'World');
Result: Hello World
```

```
SELECT format('Testing %s, %s, %s, %%', 'one', 'two', 'three');
Result: Testing one, two, three, %
```

```
SELECT format('INSERT INTO %I VALUES(%L)', 'Foo bar', E'O\Reilly');
Result: INSERT INTO "Foo bar" VALUES('O'Reilly')
```

```
SELECT format('INSERT INTO %I VALUES(%L)', 'locations', 'C:\Program Files');
Result: INSERT INTO locations VALUES('C:\Program Files')
```

Here are examples using *width* fields and the - flag:

```
SELECT format('|%10s|', 'foo');
Result: |          foo|
```

```
SELECT format('|%-10s|', 'foo');
Result: |foo          |
```

```
SELECT format('|%*s|', 10, 'foo');
Result: |          foo|
```

```
SELECT format('|%*s|', -10, 'foo');
Result: |foo          |
```

```
SELECT format('|%-*s|', 10, 'foo');
Result: |foo      |
```

```
SELECT format('|%-*s|', -10, 'foo');
Result: |foo      |
```

These examples show use of *position* fields:

```
SELECT format('Testing %3$s, %2$s, %1$s', 'one', 'two', 'three');
Result: Testing three, two, one
```

```
SELECT format('|%*2$s|', 'foo', 10, 'bar');
Result: |          bar|
```

```
SELECT format('|%1$*2$s|', 'foo', 10, 'bar');
Result: |          foo|
```

Unlike the standard C function `sprintf`, Postgres Pro's `format` function allows format specifiers with and without *position* fields to be mixed in the same format string. A format specifier without a *position* field always uses the next argument after the last argument consumed. In addition, the `format` function does not require all function arguments to be used in the format string. For example:

```
SELECT format('Testing %3$s, %2$s, %s', 'one', 'two', 'three');
Result: Testing three, two, three
```

The `%I` and `%L` format specifiers are particularly useful for safely constructing dynamic SQL statements. See [Example 44.1](#).

9.5. Binary String Functions and Operators

This section describes functions and operators for examining and manipulating binary strings, that is values of type `bytea`. Many of these are equivalent, in purpose and syntax, to the text-string functions described in the previous section.

SQL defines some string functions that use key words, rather than commas, to separate arguments. Details are in [Table 9.11](#). Postgres Pro also provides versions of these functions that use the regular function invocation syntax (see [Table 9.12](#)).

Table 9.11. SQL Binary String Functions and Operators

Function/Operator	Description	Example(s)
<code>bytea bytea</code>	<code>→ bytea</code> Concatenates the two binary strings.	<code>'\x123456'::bytea '\x789a00bcde'::bytea</code> <code>→ \x123456789a00bcde</code>
<code>bit_length (bytea)</code>	<code>→ integer</code> Returns number of bits in the binary string (8 times the <code>octet_length</code>).	<code>bit_length('\x123456'::bytea)</code> <code>→ 24</code>
<code>btrim (bytes bytea, bytesremoved bytea)</code>	<code>→ bytea</code> Removes the longest string containing only bytes appearing in <i>bytesremoved</i> from the start and end of <i>bytes</i> .	<code>btrim('\x1234567890'::bytea, '\x9012'::bytea)</code> <code>→ \x345678</code>
<code>ltrim (bytes bytea, bytesremoved bytea)</code>	<code>→ bytea</code> Removes the longest string containing only bytes appearing in <i>bytesremoved</i> from the start of <i>bytes</i> .	

Function/Operator	Description	Example(s)
		<code>ltrim('\x1234567890'::bytea, '\x9012'::bytea) → \x34567890</code>
<code>octet_length (bytea) → integer</code>	Returns number of bytes in the binary string.	<code>octet_length('\x123456'::bytea) → 3</code>
<code>overlay (bytes bytea PLACING newsubstring bytea FROM start integer [FOR count integer]) → bytea</code>	Replaces the substring of <i>bytes</i> that starts at the <i>start</i> 'th byte and extends for <i>count</i> bytes with <i>newsubstring</i> . If <i>count</i> is omitted, it defaults to the length of <i>newsubstring</i> .	<code>overlay('\x1234567890'::bytea placing '\002\003'::bytea from 2 for 3) → \x12020390</code>
<code>position (substring bytea IN bytes bytea) → integer</code>	Returns first starting index of the specified <i>substring</i> within <i>bytes</i> , or zero if it's not present.	<code>position('\x5678'::bytea in '\x1234567890'::bytea) → 3</code>
<code>rtrim (bytes bytea, bytesremoved bytea) → bytea</code>	Removes the longest string containing only bytes appearing in <i>bytesremoved</i> from the end of <i>bytes</i> .	<code>rtrim('\x1234567890'::bytea, '\x9012'::bytea) → \x12345678</code>
<code>substring (bytes bytea [FROM start integer] [FOR count integer]) → bytea</code>	Extracts the substring of <i>bytes</i> starting at the <i>start</i> 'th byte if that is specified, and stopping after <i>count</i> bytes if that is specified. Provide at least one of <i>start</i> and <i>count</i> .	<code>substring('\x1234567890'::bytea from 3 for 2) → \x5678</code>
<code>trim ([LEADING TRAILING BOTH] bytesremoved bytea FROM bytes bytea) → bytea</code>	Removes the longest string containing only bytes appearing in <i>bytesremoved</i> from the start, end, or both ends (BOTH is the default) of <i>bytes</i> .	<code>trim('\x9012'::bytea from '\x1234567890'::bytea) → \x345678</code>
<code>trim ([LEADING TRAILING BOTH] [FROM] bytes bytea, bytesremoved bytea) → bytea</code>	This is a non-standard syntax for <code>trim()</code> .	<code>trim(both from '\x1234567890'::bytea, '\x9012'::bytea) → \x345678</code>

Additional binary string manipulation functions are available and are listed in [Table 9.12](#). Some of them are used internally to implement the SQL-standard string functions listed in [Table 9.11](#).

Table 9.12. Other Binary String Functions

Function	Description	Example(s)
<code>bit_count (bytes bytea) → bigint</code>	Returns the number of bits set in the binary string (also known as “popcount”).	<code>bit_count('\x1234567890'::bytea) → 15</code>
<code>get_bit (bytes bytea, n bigint) → integer</code>	Extracts <i>n</i> 'th bit from binary string.	<code>get_bit('\x1234567890'::bytea, 30) → 1</code>
<code>get_byte (bytes bytea, n integer) → integer</code>	Extracts <i>n</i> 'th byte from binary string.	<code>get_byte('\x1234567890'::bytea, 4) → 144</code>

Function	Description	Example(s)
<code>length (bytea) → integer</code>	Returns the number of bytes in the binary string.	<code>length ('\\x1234567890'::bytea) → 5</code>
<code>length (bytes bytea, encoding name) → integer</code>	Returns the number of characters in the binary string, assuming that it is text in the given <i>encoding</i> .	<code>length ('jose'::bytea, 'UTF8') → 4</code>
<code>md5 (bytea) → text</code>	Computes the MD5 hash of the binary string, with the result written in hexadecimal.	<code>md5 ('Th\\000omas'::bytea) → 8ab2d3c9689aaf18b4958c334c82d8b1</code>
<code>set_bit (bytes bytea, n bigint, newvalue integer) → bytea</code>	Sets n th bit in binary string to <i>newvalue</i> .	<code>set_bit ('\\x1234567890'::bytea, 30, 0) → \\x1234563890</code>
<code>set_byte (bytes bytea, n integer, newvalue integer) → bytea</code>	Sets n th byte in binary string to <i>newvalue</i> .	<code>set_byte ('\\x1234567890'::bytea, 4, 64) → \\x1234567840</code>
<code>sha224 (bytea) → bytea</code>	Computes the SHA-224 hash of the binary string.	<code>sha224 ('abc'::bytea) → \\x23097d223405d8228642a477bda255b32aadbce4b- da0b3f7e36c9da7</code>
<code>sha256 (bytea) → bytea</code>	Computes the SHA-256 hash of the binary string.	<code>sha256 ('abc'::bytea) → \\xba7816bf8f01cfea414140de5dae2223 b00361a396177a9cb410ff61f20015ad</code>
<code>sha384 (bytea) → bytea</code>	Computes the SHA-384 hash of the binary string.	<code>sha384 ('abc'::bytea) → \\xcb00753f45a35e8bb5a03d699ac65007272c32ab0ed- ed1631a8b605a43ff5bed8086072ba1e7cc2358baeca134c825a7</code>
<code>sha512 (bytea) → bytea</code>	Computes the SHA-512 hash of the binary string.	<code>sha512 ('abc'::bytea) → \\xddaf35a193617abacc417349ae204131 12e6fa4e89a97ea20a9eeee64b55d39a2192992a274fc1a836ba3c23a3feebbd 454d4423643ce80e2a9ac94fa54ca49f</code>
<code>substr (bytes bytea, start integer [, count integer]) → bytea</code>	Extracts the substring of <i>bytes</i> starting at the <i>start</i> th byte, and extending for <i>count</i> bytes if that is specified. (Same as <code>substring(bytes from start for count)</code> .)	<code>substr ('\\x1234567890'::bytea, 3, 2) → \\x5678</code>

Functions `get_byte` and `set_byte` number the first byte of a binary string as byte 0. Functions `get_bit` and `set_bit` number bits from the right within each byte; for example bit 0 is the least significant bit of the first byte, and bit 15 is the most significant bit of the second byte.

For historical reasons, the function `md5` returns a hex-encoded value of type `text` whereas the SHA-2 functions return type `bytea`. Use the functions [encode](#) and [decode](#) to convert between the two. For example write `encode(sha256('abc'), 'hex')` to get a hex-encoded text representation, or `decode(md5('abc'), 'hex')` to get a `bytea` value.

Functions for converting strings between different character sets (encodings), and for representing arbitrary binary data in textual form, are shown in [Table 9.13](#). For these functions, an argument or result of type `text` is expressed in the database's default encoding, while arguments or results of type `bytea` are in an encoding named by another argument.

Table 9.13. Text/Binary String Conversion Functions

Function	Description	Example(s)
<code>convert (bytes bytea, src_encoding name, dest_encoding name) → bytea</code>	Converts a binary string representing text in encoding <i>src_encoding</i> to a binary string in encoding <i>dest_encoding</i> (see Section 23.3.4 for available conversions).	<code>convert('text_in_utf8', 'UTF8', 'LATIN1') → \x746578745f696e5f75746638</code>
<code>convert_from (bytes bytea, src_encoding name) → text</code>	Converts a binary string representing text in encoding <i>src_encoding</i> to text in the database encoding (see Section 23.3.4 for available conversions).	<code>convert_from('text_in_utf8', 'UTF8') → text_in_utf8</code>
<code>convert_to (string text, dest_encoding name) → bytea</code>	Converts a text string (in the database encoding) to a binary string encoded in encoding <i>dest_encoding</i> (see Section 23.3.4 for available conversions).	<code>convert_to('some_text', 'UTF8') → \x736f6d655f74657874</code>
<code>encode (bytes bytea, format text) → text</code>	Encodes binary data into a textual representation; supported <i>format</i> values are: base64 , escape , hex .	<code>encode('123\000\001', 'base64') → MTIzAAE=</code>
<code>decode (string text, format text) → bytea</code>	Decodes binary data from a textual representation; supported <i>format</i> values are the same as for <code>encode</code> .	<code>decode('MTIzAAE=', 'base64') → \x3132330001</code>

The `encode` and `decode` functions support the following textual formats:

base64

The `base64` format is that of [RFC 2045 Section 6.8](#). As per the RFC, encoded lines are broken at 76 characters. However instead of the MIME CRLF end-of-line marker, only a newline is used for end-of-line. The `decode` function ignores carriage-return, newline, space, and tab characters. Otherwise, an error is raised when `decode` is supplied invalid base64 data — including when trailing padding is incorrect.

escape

The `escape` format converts zero bytes and bytes with the high bit set into octal escape sequences (`\nnn`), and it doubles backslashes. Other byte values are represented literally. The `decode` function will raise an error if a backslash is not followed by either a second backslash or three octal digits; it accepts other byte values unchanged.

hex

The `hex` format represents each 4 bits of data as one hexadecimal digit, 0 through `f`, writing the higher-order digit of each byte first. The `encode` function outputs the `a-f` hex digits in lower case. Because the smallest unit of data is 8 bits, there are always an even number of characters returned by `encode`. The `decode` function accepts the `a-f` characters in either upper or lower case. An error is raised when `decode` is given invalid hex data — including when given an odd number of characters.

See also the aggregate function `string_agg` in [Section 9.21](#) and the large object functions in [Section 36.4](#).

9.6. Bit String Functions and Operators

This section describes functions and operators for examining and manipulating bit strings, that is values of the types `bit` and `bit varying`. (While only type `bit` is mentioned in these tables, values of type `bit varying` can be used interchangeably.) Bit strings support the usual comparison operators shown in [Table 9.1](#), as well as the operators shown in [Table 9.14](#).

Table 9.14. Bit String Operators

Operator	Description	Example(s)
<code>bit bit → bit</code>	Concatenation	<code>B'10001' B'011' → 10001011</code>
<code>bit & bit → bit</code>	Bitwise AND (inputs must be of equal length)	<code>B'10001' & B'01101' → 00001</code>
<code>bit bit → bit</code>	Bitwise OR (inputs must be of equal length)	<code>B'10001' B'01101' → 11101</code>
<code>bit # bit → bit</code>	Bitwise exclusive OR (inputs must be of equal length)	<code>B'10001' # B'01101' → 11100</code>
<code>~ bit → bit</code>	Bitwise NOT	<code>~ B'10001' → 01110</code>
<code>bit << integer → bit</code>	Bitwise shift left (string length is preserved)	<code>B'10001' << 3 → 01000</code>
<code>bit >> integer → bit</code>	Bitwise shift right (string length is preserved)	<code>B'10001' >> 2 → 00100</code>

Some of the functions available for binary strings are also available for bit strings, as shown in [Table 9.15](#).

Table 9.15. Bit String Functions

Function	Description	Example(s)
<code>bit_count (bit) → bigint</code>	Returns the number of bits set in the bit string (also known as “popcount”).	<code>bit_count(B'10111') → 4</code>
<code>bit_length (bit) → integer</code>	Returns number of bits in the bit string.	<code>bit_length(B'10111') → 5</code>
<code>length (bit) → integer</code>		

Function	Description	Example(s)
	Returns number of bits in the bit string.	<code>length(B'10111')</code> → 5
<code>octet_length (bit)</code>	→ integer Returns number of bytes in the bit string.	<code>octet_length(B'1011111011')</code> → 2
<code>overlay (bits bit PLACING newsubstring bit FROM start integer [FOR count integer])</code>	→ bit Replaces the substring of <i>bits</i> that starts at the <i>start</i> 'th bit and extends for <i>count</i> bits with <i>newsubstring</i> . If <i>count</i> is omitted, it defaults to the length of <i>newsubstring</i> .	<code>overlay(B'010101010101010' placing B'11111' from 2 for 3)</code> → 01111101010101010
<code>position (substring bit IN bits bit)</code>	→ integer Returns first starting index of the specified <i>substring</i> within <i>bits</i> , or zero if it's not present.	<code>position(B'010' in B'000001101011')</code> → 8
<code>substring (bits bit [FROM start integer] [FOR count integer])</code>	→ bit Extracts the substring of <i>bits</i> starting at the <i>start</i> 'th bit if that is specified, and stopping after <i>count</i> bits if that is specified. Provide at least one of <i>start</i> and <i>count</i> .	<code>substring(B'110010111111' from 3 for 2)</code> → 00
<code>get_bit (bits bit, n integer)</code>	→ integer Extracts <i>n</i> 'th bit from bit string; the first (leftmost) bit is bit 0.	<code>get_bit(B'101010101010101010', 6)</code> → 1
<code>set_bit (bits bit, n integer, newvalue integer)</code>	→ bit Sets <i>n</i> 'th bit in bit string to <i>newvalue</i> ; the first (leftmost) bit is bit 0.	<code>set_bit(B'101010101010101010', 6, 0)</code> → 101010001010101010

In addition, it is possible to cast integral values to and from type `bit`. Casting an integer to `bit(n)` copies the rightmost *n* bits. Casting an integer to a bit string width wider than the integer itself will sign-extend on the left. Some examples:

```
44::bit(10)           0000101100
44::bit(3)            100
cast(-44 as bit(12))  111111010100
'1110'::bit(4)::integer 14
```

Note that casting to just “bit” means casting to `bit(1)`, and so will deliver only the least significant bit of the integer.

9.7. Pattern Matching

There are three separate approaches to pattern matching provided by Postgres Pro: the traditional SQL `LIKE` operator, the more recent `SIMILAR TO` operator (added in SQL:1999), and POSIX-style regular expressions. Aside from the basic “does this string match this pattern?” operators, functions are available to extract or replace matching substrings and to split a string at matching locations.

Tip

If you have pattern matching needs that go beyond this, consider writing a user-defined function in Perl or Tcl.

Caution

While most regular-expression searches can be executed very quickly, regular expressions can be contrived that take arbitrary amounts of time and memory to process. Be wary of accepting regular-expression search patterns from hostile sources. If you must do so, it is advisable to impose a statement timeout.

Searches using `SIMILAR TO` patterns have the same security hazards, since `SIMILAR TO` provides many of the same capabilities as POSIX-style regular expressions.

`LIKE` searches, being much simpler than the other two options, are safer to use with possibly-hostile pattern sources.

The pattern matching operators of all three kinds do not support nondeterministic collations. If required, apply a different collation to the expression to work around this limitation.

9.7.1. LIKE

```
string LIKE pattern [ESCAPE escape-character]
string NOT LIKE pattern [ESCAPE escape-character]
```

The `LIKE` expression returns true if the *string* matches the supplied *pattern*. (As expected, the `NOT LIKE` expression returns false if `LIKE` returns true, and vice versa. An equivalent expression is `NOT (string LIKE pattern)`.)

If *pattern* does not contain percent signs or underscores, then the pattern only represents the string itself; in that case `LIKE` acts like the equals operator. An underscore (`_`) in *pattern* stands for (matches) any single character; a percent sign (`%`) matches any sequence of zero or more characters.

Some examples:

```
'abc' LIKE 'abc'      true
'abc' LIKE 'a%'       true
'abc' LIKE '_b_'      true
'abc' LIKE 'c'        false
```

`LIKE` pattern matching always covers the entire string. Therefore, if it's desired to match a sequence anywhere within a string, the pattern must start and end with a percent sign.

To match a literal underscore or percent sign without matching other characters, the respective character in *pattern* must be preceded by the escape character. The default escape character is the backslash but a different one can be selected by using the `ESCAPE` clause. To match the escape character itself, write two escape characters.

Note

If you have [standard_conforming_strings](#) turned off, any backslashes you write in literal string constants will need to be doubled. See [Section 4.1.2.1](#) for more information.

It's also possible to select no escape character by writing `ESCAPE ''`. This effectively disables the escape mechanism, which makes it impossible to turn off the special meaning of underscore and percent signs in the pattern.

According to the SQL standard, omitting `ESCAPE` means there is no escape character (rather than defaulting to a backslash), and a zero-length `ESCAPE` value is disallowed. Postgres Pro's behavior in this regard is therefore slightly nonstandard.

The key word `ILIKE` can be used instead of `LIKE` to make the match case-insensitive according to the active locale. This is not in the SQL standard but is a Postgres Pro extension.

The operator `~~` is equivalent to `LIKE`, and `~~*` corresponds to `ILIKE`. There are also `!~~` and `!~~*` operators that represent `NOT LIKE` and `NOT ILIKE`, respectively. All of these operators are Postgres Pro-specific. You may see these operator names in `EXPLAIN` output and similar places, since the parser actually translates `LIKE` et al. to these operators.

The phrases `LIKE`, `ILIKE`, `NOT LIKE`, and `NOT ILIKE` are generally treated as operators in Postgres Pro syntax; for example they can be used in *expression operator ANY (subquery)* constructs, although an `ESCAPE` clause cannot be included there. In some obscure cases it may be necessary to use the underlying operator names instead.

Also see the starts-with operator `^@` and the corresponding `starts_with()` function, which are useful in cases where simply matching the beginning of a string is needed.

9.7.2. SIMILAR TO Regular Expressions

```
string SIMILAR TO pattern [ESCAPE escape-character]
string NOT SIMILAR TO pattern [ESCAPE escape-character]
```

The `SIMILAR TO` operator returns true or false depending on whether its pattern matches the given string. It is similar to `LIKE`, except that it interprets the pattern using the SQL standard's definition of a regular expression. SQL regular expressions are a curious cross between `LIKE` notation and common (POSIX) regular expression notation.

Like `LIKE`, the `SIMILAR TO` operator succeeds only if its pattern matches the entire string; this is unlike common regular expression behavior where the pattern can match any part of the string. Also like `LIKE`, `SIMILAR TO` uses `_` and `%` as wildcard characters denoting any single character and any string, respectively (these are comparable to `.` and `.*` in POSIX regular expressions).

In addition to these facilities borrowed from `LIKE`, `SIMILAR TO` supports these pattern-matching metacharacters borrowed from POSIX regular expressions:

- `|` denotes alternation (either of two alternatives).
- `*` denotes repetition of the previous item zero or more times.
- `+` denotes repetition of the previous item one or more times.
- `?` denotes repetition of the previous item zero or one time.
- `{m}` denotes repetition of the previous item exactly *m* times.
- `{m, }` denotes repetition of the previous item *m* or more times.
- `{m, n}` denotes repetition of the previous item at least *m* and not more than *n* times.
- Parentheses `()` can be used to group items into a single logical item.
- A bracket expression `[...]` specifies a character class, just as in POSIX regular expressions.

Notice that the period `(.)` is not a metacharacter for `SIMILAR TO`.

As with `LIKE`, a backslash disables the special meaning of any of these metacharacters. A different escape character can be specified with `ESCAPE`, or the escape capability can be disabled by writing `ESCAPE ''`.

According to the SQL standard, omitting `ESCAPE` means there is no escape character (rather than defaulting to a backslash), and a zero-length `ESCAPE` value is disallowed. Postgres Pro's behavior in this regard is therefore slightly nonstandard.

Another nonstandard extension is that following the escape character with a letter or digit provides access to the escape sequences defined for POSIX regular expressions; see [Table 9.20](#), [Table 9.21](#), and [Table 9.22](#) below.

Some examples:

<code>'abc' SIMILAR TO 'abc'</code>	<code>true</code>
<code>'abc' SIMILAR TO 'a'</code>	<code>false</code>
<code>'abc' SIMILAR TO '%(b d)%'</code>	<code>true</code>
<code>'abc' SIMILAR TO '(b c)%'</code>	<code>false</code>

```
'-abc-' SIMILAR TO '%\mabc\M%' true
'xabcy' SIMILAR TO '%\mabc\M%' false
```

The `substring` function with three parameters provides extraction of a substring that matches an SQL regular expression pattern. The function can be written according to standard SQL syntax:

```
substring(string similar pattern escape escape-character)
```

or using the now obsolete SQL:1999 syntax:

```
substring(string from pattern for escape-character)
```

or as a plain three-argument function:

```
substring(string, pattern, escape-character)
```

As with `SIMILAR TO`, the specified pattern must match the entire data string, or else the function fails and returns null. To indicate the part of the pattern for which the matching data sub-string is of interest, the pattern should contain two occurrences of the escape character followed by a double quote ("). The text matching the portion of the pattern between these separators is returned when the match is successful.

The escape-double-quote separators actually divide `substring`'s pattern into three independent regular expressions; for example, a vertical bar (|) in any of the three sections affects only that section. Also, the first and third of these regular expressions are defined to match the smallest possible amount of text, not the largest, when there is any ambiguity about how much of the data string matches which pattern. (In POSIX parlance, the first and third regular expressions are forced to be non-greedy.)

As an extension to the SQL standard, Postgres Pro allows there to be just one escape-double-quote separator, in which case the third regular expression is taken as empty; or no separators, in which case the first and third regular expressions are taken as empty.

Some examples, with `#` delimiting the return string:

```
substring('foobar' similar '%#o_b#%' escape '#') oob
substring('foobar' similar '#o_b#%' escape '#') NULL
```

9.7.3. POSIX Regular Expressions

[Table 9.16](#) lists the available operators for pattern matching using POSIX regular expressions.

Table 9.16. Regular Expression Match Operators

Operator	Description	Example(s)
<code>text ~ text</code>	<code>→ boolean</code> String matches regular expression, case sensitively	<code>'thomas' ~ 't.*ma' → t</code>
<code>text ~* text</code>	<code>→ boolean</code> String matches regular expression, case-insensitively	<code>'thomas' ~* 'T.*ma' → t</code>
<code>text !~ text</code>	<code>→ boolean</code> String does not match regular expression, case sensitively	<code>'thomas' !~ 't.*max' → t</code>
<code>text !~* text</code>	<code>→ boolean</code> String does not match regular expression, case-insensitively	<code>'thomas' !~* 'T.*ma' → f</code>

POSIX regular expressions provide a more powerful means for pattern matching than the `LIKE` and `SIMILAR TO` operators. Many Unix tools such as `egrep`, `sed`, or `awk` use a pattern matching language that is similar to the one described here.

A regular expression is a character sequence that is an abbreviated definition of a set of strings (a *regular set*). A string is said to match a regular expression if it is a member of the regular set described by the regular expression. As with `LIKE`, pattern characters match string characters exactly unless they are special characters in the regular expression language — but regular expressions use different special characters than `LIKE` does. Unlike `LIKE` patterns, a regular expression is allowed to match anywhere within a string, unless the regular expression is explicitly anchored to the beginning or end of the string.

Some examples:

```
'abcd' ~ 'bc'           true
'abcd' ~ 'a.c'          true — dot matches any character
'abcd' ~ 'a.*d'         true — * repeats the preceding pattern item
'abcd' ~ '(b|x)'        true — | means OR, parentheses group
'abcd' ~ '^a'           true — ^ anchors to start of string
'abcd' ~ '^ (b|c)'      false — would match except for anchoring
```

The POSIX pattern language is described in much greater detail below.

The `substring` function with two parameters, `substring(string from pattern)`, provides extraction of a substring that matches a POSIX regular expression pattern. It returns null if there is no match, otherwise the first portion of the text that matched the pattern. But if the pattern contains any parentheses, the portion of the text that matched the first parenthesized subexpression (the one whose left parenthesis comes first) is returned. You can put parentheses around the whole expression if you want to use parentheses within it without triggering this exception. If you need parentheses in the pattern before the subexpression you want to extract, see the non-capturing parentheses described below.

Some examples:

```
substring('foobar' from 'o.b')      oob
substring('foobar' from 'o(.)b')    o
```

The `regexp_count` function counts the number of places where a POSIX regular expression pattern matches a string. It has the syntax `regexp_count(string, pattern [, start [, flags]])`. `pattern` is searched for in `string`, normally from the beginning of the string, but if the `start` parameter is provided then beginning from that character index. The `flags` parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. For example, including `i` in `flags` specifies case-insensitive matching. Supported flags are described in [Table 9.24](#).

Some examples:

```
regexp_count('ABCABCAXYaxy', 'A.')      3
regexp_count('ABCABCAXYaxy', 'A.', 1, 'i') 4
```

The `regexp_instr` function returns the starting or ending position of the *N*th match of a POSIX regular expression pattern to a string, or zero if there is no such match. It has the syntax `regexp_instr(string, pattern [, start [, N [, endoption [, flags [, subexpr]]])`. `pattern` is searched for in `string`, normally from the beginning of the string, but if the `start` parameter is provided then beginning from that character index. If *N* is specified then the *N*th match of the pattern is located, otherwise the first match is located. If the `endoption` parameter is omitted or specified as zero, the function returns the position of the first character of the match. Otherwise, `endoption` must be one, and the function returns the position of the character following the match. The `flags` parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. Supported flags are described in [Table 9.24](#). For a pattern containing parenthesized subexpressions, `subexpr` is an integer indicating which subexpression is of interest: the result identifies the position of the substring matching that subexpression. Subexpressions are numbered in the order of their leading parentheses. When `subexpr` is omitted or zero, the result identifies the position of the whole match regardless of parenthesized subexpressions.

Some examples:

```
regexp_instr('number of your street, town zip, FR', '[^,]+', 1, 2)
                                     23
regexp_instr('ABCDEFGHI', '(c..)(...)', 1, 1, 0, 'i', 2)
```

The `regexp_like` function checks whether a match of a POSIX regular expression pattern occurs within a string, returning boolean `true` or `false`. It has the syntax `regexp_like(string, pattern [, flags])`. The `flags` parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. Supported flags are described in [Table 9.24](#). This function has the same results as the `~` operator if no flags are specified. If only the `i` flag is specified, it has the same results as the `~*` operator.

Some examples:

```
regexp_like('Hello World', 'world')           false
regexp_like('Hello World', 'world', 'i')       true
```

The `regexp_match` function returns a text array of matching substring(s) within the first match of a POSIX regular expression pattern to a string. It has the syntax `regexp_match(string, pattern [, flags])`. If there is no match, the result is `NULL`. If a match is found, and the `pattern` contains no parenthesized subexpressions, then the result is a single-element text array containing the substring matching the whole pattern. If a match is found, and the `pattern` contains parenthesized subexpressions, then the result is a text array whose *n*'th element is the substring matching the *n*'th parenthesized subexpression of the `pattern` (not counting “non-capturing” parentheses; see below for details). The `flags` parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. Supported flags are described in [Table 9.24](#).

Some examples:

```
SELECT regexp_match('foobarbequebaz', 'bar.*que');
 regexp_match
-----
 {barbeque}
(1 row)

SELECT regexp_match('foobarbequebaz', '(bar) (beque)');
 regexp_match
-----
 {bar,beque}
(1 row)
```

Tip

In the common case where you just want the whole matching substring or `NULL` for no match, the best solution is to use `regexp_substr()`. However, `regexp_substr()` only exists in Postgres Pro version 15 and up. When working in older versions, you can extract the first element of `regexp_match()`'s result, for example:

```
SELECT (regexp_match('foobarbequebaz', 'bar.*que'))[1];
 regexp_match
-----
 barbeque
(1 row)
```

The `regexp_matches` function returns a set of text arrays of matching substring(s) within matches of a POSIX regular expression pattern to a string. It has the same syntax as `regexp_match`. This function returns no rows if there is no match, one row if there is a match and the `g` flag is not given, or *N* rows if there are *N* matches and the `g` flag is given. Each returned row is a text array containing the whole matched substring or the substrings matching parenthesized subexpressions of the `pattern`, just as described above for `regexp_match`. `regexp_matches` accepts all the flags shown in [Table 9.24](#), plus the `g` flag which commands it to return all matches, not just the first one.

Some examples:

```
SELECT regexp_matches('foo', 'not there');
 regexp_matches
-----
(0 rows)

SELECT regexp_matches('foobarbequebazilbarfbonk', '(b[^b]+) (b[^b]+)', 'g');
 regexp_matches
-----
{bar,beque}
{bazil,barf}
(2 rows)
```

Tip

In most cases `regexp_matches()` should be used with the `g` flag, since if you only want the first match, it's easier and more efficient to use `regexp_match()`. However, `regexp_match()` only exists in Postgres Pro version 10 and up. When working in older versions, a common trick is to place a `regexp_matches()` call in a sub-select, for example:

```
SELECT col1, (SELECT regexp_matches(col2, '(bar)(beque)')) FROM tab;
```

This produces a text array if there's a match, or `NULL` if not, the same as `regexp_match()` would do. Without the sub-select, this query would produce no output at all for table rows without a match, which is typically not the desired behavior.

The `regexp_replace` function provides substitution of new text for substrings that match POSIX regular expression patterns. It has the syntax `regexp_replace(source, pattern, replacement [, start [, N]] [, flags])`. (Notice that *N* cannot be specified unless *start* is, but *flags* can be given in any case.) The *source* string is returned unchanged if there is no match to the *pattern*. If there is a match, the *source* string is returned with the *replacement* string substituted for the matching substring. The *replacement* string can contain `\n`, where *n* is 1 through 9, to indicate that the source substring matching the *n*'th parenthesized subexpression of the pattern should be inserted, and it can contain `&` to indicate that the substring matching the entire pattern should be inserted. Write `\\` if you need to put a literal backslash in the replacement text. *pattern* is searched for in *string*, normally from the beginning of the string, but if the *start* parameter is provided then beginning from that character index. By default, only the first match of the pattern is replaced. If *N* is specified and is greater than zero, then the *N*'th match of the pattern is replaced. If the `g` flag is given, or if *N* is specified and is zero, then all matches at or after the *start* position are replaced. (The `g` flag is ignored when *N* is specified.) The *flags* parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. Supported flags (though not `g`) are described in [Table 9.24](#).

Some examples:

```
regexp_replace('foobarbaz', 'b..', 'X')
      fooXbaz
regexp_replace('foobarbaz', 'b..', 'X', 'g')
      fooXX
regexp_replace('foobarbaz', 'b(..)', 'X\\1Y', 'g')
      fooXarYXazY
regexp_replace('A PostgreSQL function', 'a|e|i|o|u', 'X', 1, 0, 'i')
      X PXstgrXSQL fXnctXXn
regexp_replace('A PostgreSQL function', 'a|e|i|o|u', 'X', 1, 3, 'i')
      A PostgrXSQL function
```

The `regexp_split_to_table` function splits a string using a POSIX regular expression pattern as a delimiter. It has the syntax `regexp_split_to_table(string, pattern [, flags])`. If there is no match to the *pattern*, the function returns the *string*. If there is at least one match, for each match it returns the text from the end of the last match (or the beginning of the string) to the beginning of the match.

When there are no more matches, it returns the text from the end of the last match to the end of the string. The *flags* parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. `regexp_split_to_table` supports the flags described in [Table 9.24](#).

The `regexp_split_to_array` function behaves the same as `regexp_split_to_table`, except that `regexp_split_to_array` returns its result as an array of text. It has the syntax `regexp_split_to_array(string, pattern [, flags])`. The parameters are the same as for `regexp_split_to_table`.

Some examples:

```
SELECT foo FROM regexp_split_to_table('the quick brown fox jumps over the lazy dog',
  '\s+') AS foo;
```

```
-----
the
quick
brown
fox
jumps
over
the
lazy
dog
(9 rows)
```

```
SELECT regexp_split_to_array('the quick brown fox jumps over the lazy dog', '\s+');
           regexp_split_to_array
```

```
-----
{the,quick,brown,fox,jumps,over,the,lazy,dog}
(1 row)
```

```
SELECT foo FROM regexp_split_to_table('the quick brown fox', '\s*') AS foo;
```

```
-----
t
h
e
q
u
i
c
k
b
r
o
w
n
f
o
x
(16 rows)
```

As the last example demonstrates, the `regexp` split functions ignore zero-length matches that occur at the start or end of the string or immediately after a previous match. This is contrary to the strict definition of `regexp` matching that is implemented by the other `regexp` functions, but is usually the most convenient behavior in practice. Other software systems such as Perl use similar definitions.

The `regexp_substr` function returns the substring that matches a POSIX regular expression pattern, or NULL if there is no match. It has the syntax `regexp_substr(string, pattern [, start [, N [, flags`

[, *subexpr*]]]). *pattern* is searched for in *string*, normally from the beginning of the string, but if the *start* parameter is provided then beginning from that character index. If *N* is specified then the *N*'th match of the pattern is returned, otherwise the first match is returned. The *flags* parameter is an optional text string containing zero or more single-letter flags that change the function's behavior. Supported flags are described in [Table 9.24](#). For a pattern containing parenthesized subexpressions, *subexpr* is an integer indicating which subexpression is of interest: the result is the substring matching that subexpression. Subexpressions are numbered in the order of their leading parentheses. When *subexpr* is omitted or zero, the result is the whole match regardless of parenthesized subexpressions.

Some examples:

```
regexp_substr('number of your street, town zip, FR', '^[^,]+', 1, 2)
                                town zip
regexp_substr('ABCDEFGHFI', '(c..)(...)', 1, 1, 'i', 2)
                                FGH
```

9.7.3.1. Regular Expression Details

Postgres Pro's regular expressions are implemented using a software package written by Henry Spencer. Much of the description of regular expressions below is copied verbatim from his manual.

Regular expressions (REs), as defined in POSIX 1003.2, come in two forms: *extended* REs or EREs (roughly those of `egrep`), and *basic* REs or BREs (roughly those of `ed`). Postgres Pro supports both forms, and also implements some extensions that are not in the POSIX standard, but have become widely used due to their availability in programming languages such as Perl and Tcl. REs using these non-POSIX extensions are called *advanced* REs or AREs in this documentation. AREs are almost an exact superset of EREs, but BREs have several notational incompatibilities (as well as being much more limited). We first describe the ARE and ERE forms, noting features that apply only to AREs, and then describe how BREs differ.

Note

Postgres Pro always initially presumes that a regular expression follows the ARE rules. However, the more limited ERE or BRE rules can be chosen by prepending an *embedded option* to the RE pattern, as described in [Section 9.7.3.4](#). This can be useful for compatibility with applications that expect exactly the POSIX 1003.2 rules.

A regular expression is defined as one or more *branches*, separated by `|`. It matches anything that matches one of the branches.

A branch is zero or more *quantified atoms* or *constraints*, concatenated. It matches a match for the first, followed by a match for the second, etc.; an empty branch matches the empty string.

A quantified atom is an *atom* possibly followed by a single *quantifier*. Without a quantifier, it matches a match for the atom. With a quantifier, it can match some number of matches of the atom. An *atom* can be any of the possibilities shown in [Table 9.17](#). The possible quantifiers and their meanings are shown in [Table 9.18](#).

A *constraint* matches an empty string, but matches only when specific conditions are met. A constraint can be used where an atom could be used, except it cannot be followed by a quantifier. The simple constraints are shown in [Table 9.19](#); some more constraints are described later.

Table 9.17. Regular Expression Atoms

Atom	Description
(<i>re</i>)	(where <i>re</i> is any regular expression) matches a match for <i>re</i> , with the match noted for possible reporting

Atom	Description
(?: <i>re</i>)	as above, but the match is not noted for reporting (a “non-capturing” set of parentheses) (AREs only)
.	matches any single character
[<i>chars</i>]	a <i>bracket expression</i> , matching any one of the <i>chars</i> (see Section 9.7.3.2 for more detail)
\k	(where <i>k</i> is a non-alphanumeric character) matches that character taken as an ordinary character, e.g., \\ matches a backslash character
\c	where <i>c</i> is alphanumeric (possibly followed by other characters) is an <i>escape</i> , see Section 9.7.3.3 (AREs only; in EREs and BREs, this matches <i>c</i>)
{	when followed by a character other than a digit, matches the left-brace character {; when followed by a digit, it is the beginning of a <i>bound</i> (see below)
<i>x</i>	where <i>x</i> is a single character with no other significance, matches that character

An RE cannot end with a backslash (\).

Note

If you have [standard_conforming_strings](#) turned off, any backslashes you write in literal string constants will need to be doubled. See [Section 4.1.2.1](#) for more information.

Table 9.18. Regular Expression Quantifiers

Quantifier	Matches
*	a sequence of 0 or more matches of the atom
+	a sequence of 1 or more matches of the atom
?	a sequence of 0 or 1 matches of the atom
{ <i>m</i> }	a sequence of exactly <i>m</i> matches of the atom
{ <i>m</i> , }	a sequence of <i>m</i> or more matches of the atom
{ <i>m</i> , <i>n</i> }	a sequence of <i>m</i> through <i>n</i> (inclusive) matches of the atom; <i>m</i> cannot exceed <i>n</i>
*?	non-greedy version of *
+?	non-greedy version of +
??	non-greedy version of ?
{ <i>m</i> }?	non-greedy version of { <i>m</i> }
{ <i>m</i> , }?	non-greedy version of { <i>m</i> , }
{ <i>m</i> , <i>n</i> }?	non-greedy version of { <i>m</i> , <i>n</i> }

The forms using { . . . } are known as *bounds*. The numbers *m* and *n* within a bound are unsigned decimal integers with permissible values from 0 to 255 inclusive.

Non-greedy quantifiers (available in AREs only) match the same possibilities as their corresponding normal (*greedy*) counterparts, but prefer the smallest number rather than the largest number of matches. See [Section 9.7.3.5](#) for more detail.

Note

A quantifier cannot immediately follow another quantifier, e.g., `**` is invalid. A quantifier cannot begin an expression or subexpression or follow `^` or `|`.

Table 9.19. Regular Expression Constraints

Constraint	Description
<code>^</code>	matches at the beginning of the string
<code>\$</code>	matches at the end of the string
<code>(?= re)</code>	<i>positive lookahead</i> matches at any point where a substring matching <i>re</i> begins (AREs only)
<code>(?! re)</code>	<i>negative lookahead</i> matches at any point where no substring matching <i>re</i> begins (AREs only)
<code>(?<= re)</code>	<i>positive lookbehind</i> matches at any point where a substring matching <i>re</i> ends (AREs only)
<code>(?<! re)</code>	<i>negative lookbehind</i> matches at any point where no substring matching <i>re</i> ends (AREs only)

Lookahead and lookbehind constraints cannot contain *back references* (see [Section 9.7.3.3](#)), and all parentheses within them are considered non-capturing.

9.7.3.2. Bracket Expressions

A *bracket expression* is a list of characters enclosed in `[]`. It normally matches any single character from the list (but see below). If the list begins with `^`, it matches any single character *not* from the rest of the list. If two characters in the list are separated by `-`, this is shorthand for the full range of characters between those two (inclusive) in the collating sequence, e.g., `[0-9]` in ASCII matches any decimal digit. It is illegal for two ranges to share an endpoint, e.g., `a-c-e`. Ranges are very collating-sequence-dependent, so portable programs should avoid relying on them.

To include a literal `]` in the list, make it the first character (after `^`, if that is used). To include a literal `-`, make it the first or last character, or the second endpoint of a range. To use a literal `-` as the first endpoint of a range, enclose it in `[.]` and `.]` to make it a collating element (see below). With the exception of these characters, some combinations using `[` (see next paragraphs), and escapes (AREs only), all other special characters lose their special significance within a bracket expression. In particular, `\` is not special when following ERE or BRE rules, though it is special (as introducing an escape) in AREs.

Within a bracket expression, a collating element (a character, a multiple-character sequence that collates as if it were a single character, or a collating-sequence name for either) enclosed in `[.]` and `.]` stands for the sequence of characters of that collating element. The sequence is treated as a single element of the bracket expression's list. This allows a bracket expression containing a multiple-character collating element to match more than one character, e.g., if the collating sequence includes a `ch` collating element, then the RE `[[.ch.]] *c` matches the first five characters of `chchcc`.

Note

Postgres Pro currently does not support multi-character collating elements. This information describes possible future behavior.

Within a bracket expression, a collating element enclosed in `[=` and `=]` is an *equivalence class*, standing for the sequences of characters of all collating elements equivalent to that one, including itself. (If there are no other equivalent collating elements, the treatment is as if the enclosing delimiters were `[.]` and `.]`.) For example, if `o` and `^` are the members of an equivalence class, then `[[=o=]]`, `[[=^=]]`, and `[o^]` are all synonymous. An equivalence class cannot be an endpoint of a range.

Within a bracket expression, the name of a character class enclosed in `[:` and `:]` stands for the list of all characters belonging to that class. A character class cannot be used as an endpoint of a range. The POSIX standard defines these character class names: `alnum` (letters and numeric digits), `alpha` (letters), `blank` (space and tab), `cntrl` (control characters), `digit` (numeric digits), `graph` (printable characters except space), `lower` (lower-case letters), `print` (printable characters including space), `punct` (punctuation), `space` (any white space), `upper` (upper-case letters), and `xdigit` (hexadecimal digits). The behavior of these standard character classes is generally consistent across platforms for characters in the 7-bit ASCII set. Whether a given non-ASCII character is considered to belong to one of these classes depends on the *collation* that is used for the regular-expression function or operator (see [Section 23.2](#)), or by default on the database's `LC_CTYPE` locale setting (see [Section 23.1](#)). The classification of non-ASCII characters can vary across platforms even in similarly-named locales. (But the `C` locale never considers any non-ASCII characters to belong to any of these classes.) In addition to these standard character classes, Postgres Pro defines the `word` character class, which is the same as `alnum` plus the underscore (`_`) character, and the `ascii` character class, which contains exactly the 7-bit ASCII set.

There are two special cases of bracket expressions: the bracket expressions `[[:<:]]` and `[[:>:]]` are constraints, matching empty strings at the beginning and end of a word respectively. A word is defined as a sequence of word characters that is neither preceded nor followed by word characters. A word character is any character belonging to the `word` character class, that is, any letter, digit, or underscore. This is an extension, compatible with but not specified by POSIX 1003.2, and should be used with caution in software intended to be portable to other systems. The constraint escapes described below are usually preferable; they are no more standard, but are easier to type.

9.7.3.3. Regular Expression Escapes

Escapes are special sequences beginning with `\` followed by an alphanumeric character. Escapes come in several varieties: character entry, class shorthands, constraint escapes, and back references. A `\` followed by an alphanumeric character but not constituting a valid escape is illegal in AREs. In EREs, there are no escapes: outside a bracket expression, a `\` followed by an alphanumeric character merely stands for that character as an ordinary character, and inside a bracket expression, `\` is an ordinary character. (The latter is the one actual incompatibility between EREs and AREs.)

Character-entry escapes exist to make it easier to specify non-printing and other inconvenient characters in REs. They are shown in [Table 9.20](#).

Class-shorthand escapes provide shorthands for certain commonly-used character classes. They are shown in [Table 9.21](#).

A *constraint escape* is a constraint, matching the empty string if specific conditions are met, written as an escape. They are shown in [Table 9.22](#).

A *back reference* (`\n`) matches the same string matched by the previous parenthesized subexpression specified by the number *n* (see [Table 9.23](#)). For example, `([bc])\1` matches `bb` or `cc` but not `bc` or `cb`. The subexpression must entirely precede the back reference in the RE. Subexpressions are numbered in the order of their leading parentheses. Non-capturing parentheses do not define subexpressions. The back reference considers only the string characters matched by the referenced subexpression, not any constraints contained in it. For example, `(^\d)\1` will match `22`.

Table 9.20. Regular Expression Character-Entry Escapes

Escape	Description
<code>\a</code>	alert (bell) character, as in C
<code>\b</code>	backspace, as in C
<code>\B</code>	synonym for backslash (<code>\</code>) to help reduce the need for backslash doubling
<code>\cX</code>	(where <i>x</i> is any character) the character whose low-order 5 bits are the same as those of <i>x</i> , and whose other bits are all zero

Escape	Description
<code>\e</code>	the character whose collating-sequence name is ESC, or failing that, the character with octal value 033
<code>\f</code>	form feed, as in C
<code>\n</code>	newline, as in C
<code>\r</code>	carriage return, as in C
<code>\t</code>	horizontal tab, as in C
<code>\uwxxyz</code>	(where <i>wxyz</i> is exactly four hexadecimal digits) the character whose hexadecimal value is <code>0xwxxyz</code>
<code>\Ustuvwxyz</code>	(where <i>stuvwxyz</i> is exactly eight hexadecimal digits) the character whose hexadecimal value is <code>0xstuvwxyz</code>
<code>\v</code>	vertical tab, as in C
<code>\xhhh</code>	(where <i>hhh</i> is any sequence of hexadecimal digits) the character whose hexadecimal value is <code>0xhhh</code> (a single character no matter how many hexadecimal digits are used)
<code>\0</code>	the character whose value is 0 (the null byte)
<code>\xy</code>	(where <i>xy</i> is exactly two octal digits, and is not a <i>back reference</i>) the character whose octal value is <code>0xy</code>
<code>\xyz</code>	(where <i>xyz</i> is exactly three octal digits, and is not a <i>back reference</i>) the character whose octal value is <code>0xyz</code>

Hexadecimal digits are 0-9, a-f, and A-F. Octal digits are 0-7.

Numeric character-entry escapes specifying values outside the ASCII range (0-127) have meanings dependent on the database encoding. When the encoding is UTF-8, escape values are equivalent to Unicode code points, for example `\u1234` means the character U+1234. For other multibyte encodings, character-entry escapes usually just specify the concatenation of the byte values for the character. If the escape value does not correspond to any legal character in the database encoding, no error will be raised, but it will never match any data.

The character-entry escapes are always taken as ordinary characters. For example, `\135` is `]` in ASCII, but `\135` does not terminate a bracket expression.

Table 9.21. Regular Expression Class-Shorthand Escapes

Escape	Description
<code>\d</code>	matches any digit, like <code>[[:digit:]]</code>
<code>\s</code>	matches any whitespace character, like <code>[[:space:]]</code>
<code>\w</code>	matches any word character, like <code>[[:word:]]</code>
<code>\D</code>	matches any non-digit, like <code>[^[:digit:]]</code>
<code>\S</code>	matches any non-whitespace character, like <code>[^[:space:]]</code>
<code>\W</code>	matches any non-word character, like <code>[^[:word:]]</code>

The class-shorthand escapes also work within bracket expressions, although the definitions shown above are not quite syntactically valid in that context. For example, `[a-c\d]` is equivalent to `[a-c[:digit:]]`.

Table 9.22. Regular Expression Constraint Escapes

Escape	Description
<code>\A</code>	matches only at the beginning of the string (see Section 9.7.3.5 for how this differs from <code>^</code>)
<code>\m</code>	matches only at the beginning of a word
<code>\M</code>	matches only at the end of a word
<code>\y</code>	matches only at the beginning or end of a word
<code>\Y</code>	matches only at a point that is not the beginning or end of a word
<code>\Z</code>	matches only at the end of the string (see Section 9.7.3.5 for how this differs from <code>\$</code>)

A word is defined as in the specification of `[[:<:]]` and `[[:>:]]` above. Constraint escapes are illegal within bracket expressions.

Table 9.23. Regular Expression Back References

Escape	Description
<code>\m</code>	(where <i>m</i> is a nonzero digit) a back reference to the <i>m</i> 'th subexpression
<code>\mnn</code>	(where <i>m</i> is a nonzero digit, and <i>nn</i> is some more digits, and the decimal value <i>mnn</i> is not greater than the number of closing capturing parentheses seen so far) a back reference to the <i>mnn</i> 'th subexpression

Note

There is an inherent ambiguity between octal character-entry escapes and back references, which is resolved by the following heuristics, as hinted at above. A leading zero always indicates an octal escape. A single non-zero digit, not followed by another digit, is always taken as a back reference. A multi-digit sequence not starting with a zero is taken as a back reference if it comes after a suitable subexpression (i.e., the number is in the legal range for a back reference), and otherwise is taken as octal.

9.7.3.4. Regular Expression Metasyntax

In addition to the main syntax described above, there are some special forms and miscellaneous syntactic facilities available.

An RE can begin with one of two special *director* prefixes. If an RE begins with `***:`, the rest of the RE is taken as an ARE. (This normally has no effect in Postgres Pro, since REs are assumed to be AREs; but it does have an effect if ERE or BRE mode had been specified by the *flags* parameter to a regex function.) If an RE begins with `***=`, the rest of the RE is taken to be a literal string, with all characters considered ordinary characters.

An ARE can begin with *embedded options*: a sequence `(?xyz)` (where *xyz* is one or more alphabetic characters) specifies options affecting the rest of the RE. These options override any previously determined options — in particular, they can override the case-sensitivity behavior implied by a regex operator, or the *flags* parameter to a regex function. The available option letters are shown in [Table 9.24](#). Note that these same option letters are used in the *flags* parameters of regex functions.

Table 9.24. ARE Embedded-Option Letters

Option	Description
b	rest of RE is a BRE
c	case-sensitive matching (overrides operator type)
e	rest of RE is an ERE
i	case-insensitive matching (see Section 9.7.3.5) (overrides operator type)
m	historical synonym for n
n	newline-sensitive matching (see Section 9.7.3.5)
p	partial newline-sensitive matching (see Section 9.7.3.5)
q	rest of RE is a literal (“quoted”) string, all ordinary characters
s	non-newline-sensitive matching (default)
t	tight syntax (default; see below)
w	inverse partial newline-sensitive (“weird”) matching (see Section 9.7.3.5)
x	expanded syntax (see below)

Embedded options take effect at the `)` terminating the sequence. They can appear only at the start of an ARE (after the `***:` director if any).

In addition to the usual (*tight*) RE syntax, in which all characters are significant, there is an *expanded* syntax, available by specifying the embedded `x` option. In the expanded syntax, white-space characters in the RE are ignored, as are all characters between a `#` and the following newline (or the end of the RE). This permits paragraphing and commenting a complex RE. There are three exceptions to that basic rule:

- a white-space character or `#` preceded by `\` is retained
- white space or `#` within a bracket expression is retained
- white space and comments cannot appear within multi-character symbols, such as `(?:`

For this purpose, white-space characters are blank, tab, newline, and any character that belongs to the *space* character class.

Finally, in an ARE, outside bracket expressions, the sequence `(?#ttt)` (where *ttt* is any text not containing a `)`) is a comment, completely ignored. Again, this is not allowed between the characters of multi-character symbols, like `(?:`. Such comments are more a historical artifact than a useful facility, and their use is deprecated; use the expanded syntax instead.

None of these metasyntax extensions is available if an initial `***=` director has specified that the user's input be treated as a literal string rather than as an RE.

9.7.3.5. Regular Expression Matching Rules

In the event that an RE could match more than one substring of a given string, the RE matches the one starting earliest in the string. If the RE could match more than one substring starting at that point, either the longest possible match or the shortest possible match will be taken, depending on whether the RE is *greedy* or *non-greedy*.

Whether an RE is greedy or not is determined by the following rules:

- Most atoms, and all constraints, have no greediness attribute (because they cannot match variable amounts of text anyway).
- Adding parentheses around an RE does not change its greediness.

- A quantified atom with a fixed-repetition quantifier ($\{m\}$ or $\{m\}?$) has the same greediness (possibly none) as the atom itself.
- A quantified atom with other normal quantifiers (including $\{m, n\}$ with m equal to n) is greedy (prefers longest match).
- A quantified atom with a non-greedy quantifier (including $\{m, n\}?$ with m equal to n) is non-greedy (prefers shortest match).
- A branch — that is, an RE that has no top-level `|` operator — has the same greediness as the first quantified atom in it that has a greediness attribute.
- An RE consisting of two or more branches connected by the `|` operator is always greedy.

The above rules associate greediness attributes not only with individual quantified atoms, but with branches and entire REs that contain quantified atoms. What that means is that the matching is done in such a way that the branch, or whole RE, matches the longest or shortest possible substring *as a whole*. Once the length of the entire match is determined, the part of it that matches any particular subexpression is determined on the basis of the greediness attribute of that subexpression, with subexpressions starting earlier in the RE taking priority over ones starting later.

An example of what this means:

```
SELECT SUBSTRING('XY1234Z', 'Y*([0-9]{1,3})');
Result: 123
SELECT SUBSTRING('XY1234Z', 'Y*?([0-9]{1,3})');
Result: 1
```

In the first case, the RE as a whole is greedy because `Y*` is greedy. It can match beginning at the `Y`, and it matches the longest possible string starting there, i.e., `Y123`. The output is the parenthesized part of that, or `123`. In the second case, the RE as a whole is non-greedy because `Y*?` is non-greedy. It can match beginning at the `Y`, and it matches the shortest possible string starting there, i.e., `Y1`. The subexpression `[0-9]{1,3}` is greedy but it cannot change the decision as to the overall match length; so it is forced to match just `1`.

In short, when an RE contains both greedy and non-greedy subexpressions, the total match length is either as long as possible or as short as possible, according to the attribute assigned to the whole RE. The attributes assigned to the subexpressions only affect how much of that match they are allowed to “eat” relative to each other.

The quantifiers $\{1, 1\}$ and $\{1, 1\}?$ can be used to force greediness or non-greediness, respectively, on a subexpression or a whole RE. This is useful when you need the whole RE to have a greediness attribute different from what's deduced from its elements. As an example, suppose that we are trying to separate a string containing some digits into the digits and the parts before and after them. We might try to do that like this:

```
SELECT regexp_match('abc01234xyz', '(.*) (\d+) (.*)');
Result: {abc0123,4,xyz}
```

That didn't work: the first `.*` is greedy so it “eats” as much as it can, leaving the `\d+` to match at the last possible place, the last digit. We might try to fix that by making it non-greedy:

```
SELECT regexp_match('abc01234xyz', '(.*)? (\d+) (.*)');
Result: {abc,0,""}
```

That didn't work either, because now the RE as a whole is non-greedy and so it ends the overall match as soon as possible. We can get what we want by forcing the RE as a whole to be greedy:

```
SELECT regexp_match('abc01234xyz', '(?:.*)? (\d+) (.*) {1,1}');
Result: {abc,01234,xyz}
```

Controlling the RE's overall greediness separately from its components' greediness allows great flexibility in handling variable-length patterns.

When deciding what is a longer or shorter match, match lengths are measured in characters, not collating elements. An empty string is considered longer than no match at all. For example: `bb*` matches

the three middle characters of `abbbc`; `(week|wee)(night|knights)` matches all ten characters of `week-nights`; when `(.*)` is matched against `abc` the parenthesized subexpression matches all three characters; and when `(a*)` is matched against `bc` both the whole RE and the parenthesized subexpression match an empty string.

If case-independent matching is specified, the effect is much as if all case distinctions had vanished from the alphabet. When an alphabetic that exists in multiple cases appears as an ordinary character outside a bracket expression, it is effectively transformed into a bracket expression containing both cases, e.g., `x` becomes `[xX]`. When it appears inside a bracket expression, all case counterparts of it are added to the bracket expression, e.g., `[x]` becomes `[xX]` and `[^x]` becomes `[^xX]`.

If newline-sensitive matching is specified, `.` and bracket expressions using `^` will never match the newline character (so that matches will not cross lines unless the RE explicitly includes a newline) and `^` and `$` will match the empty string after and before a newline respectively, in addition to matching at beginning and end of string respectively. But the ARE escapes `\A` and `\Z` continue to match beginning or end of string *only*. Also, the character class shorthands `\D` and `\W` will match a newline regardless of this mode. (Before Postgres Pro 14, they did not match newlines when in newline-sensitive mode. Write `[^[:digit:]]` or `[^[:word:]]` to get the old behavior.)

If partial newline-sensitive matching is specified, this affects `.` and bracket expressions as with newline-sensitive matching, but not `^` and `$`.

If inverse partial newline-sensitive matching is specified, this affects `^` and `$` as with newline-sensitive matching, but not `.` and bracket expressions. This isn't very useful but is provided for symmetry.

9.7.3.6. Limits and Compatibility

No particular limit is imposed on the length of REs in this implementation. However, programs intended to be highly portable should not employ REs longer than 256 bytes, as a POSIX-compliant implementation can refuse to accept such REs.

The only feature of AREs that is actually incompatible with POSIX EREs is that `\` does not lose its special significance inside bracket expressions. All other ARE features use syntax which is illegal or has undefined or unspecified effects in POSIX EREs; the `***` syntax of directors likewise is outside the POSIX syntax for both BREs and EREs.

Many of the ARE extensions are borrowed from Perl, but some have been changed to clean them up, and a few Perl extensions are not present. Incompatibilities of note include `\b`, `\B`, the lack of special treatment for a trailing newline, the addition of complemented bracket expressions to the things affected by newline-sensitive matching, the restrictions on parentheses and back references in lookahead/look-behind constraints, and the longest/shortest-match (rather than first-match) matching semantics.

9.7.3.7. Basic Regular Expressions

BREs differ from EREs in several respects. In BREs, `|`, `+`, and `?` are ordinary characters and there is no equivalent for their functionality. The delimiters for bounds are `\{` and `\}`, with `{` and `}` by themselves ordinary characters. The parentheses for nested subexpressions are `\(` and `\)`, with `(` and `)` by themselves ordinary characters. `^` is an ordinary character except at the beginning of the RE or the beginning of a parenthesized subexpression, `$` is an ordinary character except at the end of the RE or the end of a parenthesized subexpression, and `*` is an ordinary character if it appears at the beginning of the RE or the beginning of a parenthesized subexpression (after a possible leading `^`). Finally, single-digit back references are available, and `\<` and `\>` are synonyms for `[[:<:]]` and `[[:>:]]` respectively; no other escapes are available in BREs.

9.7.3.8. Differences from SQL Standard and XQuery

Since SQL:2008, the SQL standard includes regular expression operators and functions that performs pattern matching according to the XQuery regular expression standard:

- `LIKE_REGEX`
- `OCCURRENCES_REGEX`

- POSITION_REGEX
- SUBSTRING_REGEX
- TRANSLATE_REGEX

Postgres Pro does not currently implement these operators and functions. You can get approximately equivalent functionality in each case as shown in [Table 9.25](#). (Various optional clauses on both sides have been omitted in this table.)

Table 9.25. Regular Expression Functions Equivalencies

SQL standard	PostgreSQL
<i>string</i> LIKE_REGEX <i>pattern</i>	regexp_like(<i>string</i> , <i>pattern</i>) OR <i>string</i> ~ <i>pattern</i>
OCCURRENCES_REGEX(<i>pattern</i> IN <i>string</i>)	regexp_count(<i>string</i> , <i>pattern</i>)
POSITION_REGEX(<i>pattern</i> IN <i>string</i>)	regexp_instr(<i>string</i> , <i>pattern</i>)
SUBSTRING_REGEX(<i>pattern</i> IN <i>string</i>)	regexp_substr(<i>string</i> , <i>pattern</i>)
TRANSLATE_REGEX(<i>pattern</i> IN <i>string</i> WITH <i>replacement</i>)	regexp_replace(<i>string</i> , <i>pattern</i> , <i>re-</i> <i>placement</i>)

Regular expression functions similar to those provided by PostgreSQL are also available in a number of other SQL implementations, whereas the SQL-standard functions are not as widely implemented. Some of the details of the regular expression syntax will likely differ in each implementation.

The SQL-standard operators and functions use XQuery regular expressions, which are quite close to the ARE syntax described above. Notable differences between the existing POSIX-based regular-expression feature and XQuery regular expressions include:

- XQuery character class subtraction is not supported. An example of this feature is using the following to match only English consonants: `[a-z-[aeiou]]`.
- XQuery character class shorthands `\c`, `\C`, `\i`, and `\I` are not supported.
- XQuery character class elements using `\p{UnicodeProperty}` or the inverse `\P{UnicodeProperty}` are not supported.
- POSIX interprets character classes such as `\w` (see [Table 9.21](#)) according to the prevailing locale (which you can control by attaching a `COLLATE` clause to the operator or function). XQuery specifies these classes by reference to Unicode character properties, so equivalent behavior is obtained only with a locale that follows the Unicode rules.
- The SQL standard (not XQuery itself) attempts to cater for more variants of “newline” than POSIX does. The newline-sensitive matching options described above consider only ASCII NL (`\n`) to be a newline, but SQL would have us treat CR (`\r`), CRLF (`\r\n`) (a Windows-style newline), and some Unicode-only characters like LINE SEPARATOR (U+2028) as newlines as well. Notably, `.` and `\s` should count `\r\n` as one character not two according to SQL.
- Of the character-entry escapes described in [Table 9.20](#), XQuery supports only `\n`, `\r`, and `\t`.
- XQuery does not support the `[:name:]` syntax for character classes within bracket expressions.
- XQuery does not have lookahead or lookbehind constraints, nor any of the constraint escapes described in [Table 9.22](#).
- The metasyntax forms described in [Section 9.7.3.4](#) do not exist in XQuery.
- The regular expression flag letters defined by XQuery are related to but not the same as the option letters for POSIX ([Table 9.24](#)). While the `i` and `q` options behave the same, others do not:
 - XQuery's `s` (allow dot to match newline) and `m` (allow `^` and `$` to match at newlines) flags provide access to the same behaviors as POSIX's `n`, `p` and `w` flags, but they do *not* match the behavior of POSIX's `s` and `m` flags. Note in particular that dot-matches-newline is the default behavior in POSIX but not XQuery.

- XQuery's `x` (ignore whitespace in pattern) flag is noticeably different from POSIX's `expand-mode` flag. POSIX's `x` flag also allows `#` to begin a comment in the pattern, and POSIX will not ignore a whitespace character after a backslash.

9.8. Data Type Formatting Functions

The Postgres Pro formatting functions provide a powerful set of tools for converting various data types (date/time, integer, floating point, numeric) to formatted strings and for converting from formatted strings to specific data types. [Table 9.26](#) lists them. These functions all follow a common calling convention: the first argument is the value to be formatted and the second argument is a template that defines the output or input format.

Table 9.26. Formatting Functions

Function	Description	Example(s)
<code>to_char (timestamp, text)</code>	→ text	
<code>to_char (timestamp with time zone, text)</code>	→ text	
	Converts time stamp to string according to the given format.	
	<code>to_char(timestamp '2002-04-20 17:31:12.66', 'HH12:MI:SS')</code>	→ 05:31:12
<code>to_char (interval, text)</code>	→ text	
	Converts interval to string according to the given format.	
	<code>to_char(interval '15h 2m 12s', 'HH24:MI:SS')</code>	→ 15:02:12
<code>to_char (numeric_type , text)</code>	→ text	
	Converts number to string according to the given format; available for integer, bigint, numeric, real, double precision.	
	<code>to_char(125, '999')</code>	→ 125
	<code>to_char(125.8::real, '999D9')</code>	→ 125.8
	<code>to_char(-125.8, '999D99S')</code>	→ 125.80-
<code>to_date (text, text)</code>	→ date	
	Converts string to date according to the given format.	
	<code>to_date('05 Dec 2000', 'DD Mon YYYY')</code>	→ 2000-12-05
<code>to_number (text, text)</code>	→ numeric	
	Converts string to numeric according to the given format.	
	<code>to_number('12,454.8-', '99G999D9S')</code>	→ -12454.8
<code>to_timestamp (text, text)</code>	→ timestamp with time zone	
	Converts string to time stamp according to the given format. (See also <code>to_timestamp(double precision)</code> in Table 9.33 .)	
	<code>to_timestamp('05 Dec 2000', 'DD Mon YYYY')</code>	→ 2000-12-05 00:00:00-05

Tip

`to_timestamp` and `to_date` exist to handle input formats that cannot be converted by simple casting. For most standard date/time formats, simply casting the source string to the required data type works, and is much easier. Similarly, `to_number` is unnecessary for standard numeric representations.

In a `to_char` output template string, there are certain patterns that are recognized and replaced with appropriately-formatted data based on the given value. Any text that is not a template pattern is simply copied verbatim. Similarly, in an input template string (for the other functions), template patterns iden-

tify the values to be supplied by the input data string. If there are characters in the template string that are not template patterns, the corresponding characters in the input data string are simply skipped over (whether or not they are equal to the template string characters).

[Table 9.27](#) shows the template patterns available for formatting date and time values.

Table 9.27. Template Patterns for Date/Time Formatting

Pattern	Description
HH	hour of day (01-12)
HH12	hour of day (01-12)
HH24	hour of day (00-23)
MI	minute (00-59)
SS	second (00-59)
MS	millisecond (000-999)
US	microsecond (000000-999999)
FF1	tenth of second (0-9)
FF2	hundredth of second (00-99)
FF3	millisecond (000-999)
FF4	tenth of a millisecond (0000-9999)
FF5	hundredth of a millisecond (00000-99999)
FF6	microsecond (000000-999999)
SSSS, SSSSS	seconds past midnight (0-86399)
AM, am, PM or pm	meridiem indicator (without periods)
A.M., a.m., P.M. or p.m.	meridiem indicator (with periods)
Y, YYYY	year (4 or more digits) with comma
YYYY	year (4 or more digits)
YYY	last 3 digits of year
YY	last 2 digits of year
Y	last digit of year
IYYY	ISO 8601 week-numbering year (4 or more digits)
IYY	last 3 digits of ISO 8601 week-numbering year
IY	last 2 digits of ISO 8601 week-numbering year
I	last digit of ISO 8601 week-numbering year
BC, bc, AD or ad	era indicator (without periods)
B.C., b.c., A.D. or a.d.	era indicator (with periods)
MONTH	full upper case month name (blank-padded to 9 chars)
Month	full capitalized month name (blank-padded to 9 chars)
month	full lower case month name (blank-padded to 9 chars)
MON	abbreviated upper case month name (3 chars in English, localized lengths vary)
Mon	abbreviated capitalized month name (3 chars in English, localized lengths vary)

Pattern	Description
mon	abbreviated lower case month name (3 chars in English, localized lengths vary)
MM	month number (01-12)
DAY	full upper case day name (blank-padded to 9 chars)
Day	full capitalized day name (blank-padded to 9 chars)
day	full lower case day name (blank-padded to 9 chars)
DY	abbreviated upper case day name (3 chars in English, localized lengths vary)
Dy	abbreviated capitalized day name (3 chars in English, localized lengths vary)
dy	abbreviated lower case day name (3 chars in English, localized lengths vary)
DDD	day of year (001-366)
IDDD	day of ISO 8601 week-numbering year (001-371; day 1 of the year is Monday of the first ISO week)
DD	day of month (01-31)
D	day of the week, Sunday (1) to Saturday (7)
ID	ISO 8601 day of the week, Monday (1) to Sunday (7)
W	week of month (1-5) (the first week starts on the first day of the month)
WW	week number of year (1-53) (the first week starts on the first day of the year)
IW	week number of ISO 8601 week-numbering year (01-53; the first Thursday of the year is in week 1)
CC	century (2 digits) (the twenty-first century starts on 2001-01-01)
J	Julian Date (integer days since November 24, 4714 BC at local midnight; see Section B.7)
Q	quarter
RM	month in upper case Roman numerals (I-XII; I=January)
rm	month in lower case Roman numerals (i-xii; i=January)
TZ	upper case time-zone abbreviation
tz	lower case time-zone abbreviation
TZH	time-zone hours
TZM	time-zone minutes
OF	time-zone offset from UTC (<i>HH</i> or <i>HH:MM</i>)

Modifiers can be applied to any template pattern to alter its behavior. For example, `FMMonth` is the `Month` pattern with the `FM` modifier. [Table 9.28](#) shows the modifier patterns for date/time formatting.

Table 9.28. Template Pattern Modifiers for Date/Time Formatting

Modifier	Description	Example
FM prefix	fill mode (suppress leading zeroes and padding blanks)	FMMonth
TH suffix	upper case ordinal number suffix	DDTH, e.g., 12TH
th suffix	lower case ordinal number suffix	DDth, e.g., 12th
FX prefix	fixed format global option (see usage notes)	FX Month DD Day
TM prefix	translation mode (use localized day and month names based on lc_time)	TMMonth
SP suffix	spell mode (not implemented)	DDSP

Usage notes for date/time formatting:

- FM suppresses leading zeroes and trailing blanks that would otherwise be added to make the output of a pattern be fixed-width. In Postgres Pro, FM modifies only the next specification, while in Oracle FM affects all subsequent specifications, and repeated FM modifiers toggle fill mode on and off.
- TM suppresses trailing blanks whether or not FM is specified.
- to_timestamp and to_date ignore letter case in the input; so for example MON, Mon, and mon all accept the same strings. When using the TM modifier, case-folding is done according to the rules of the function's input collation (see [Section 23.2](#)).
- to_timestamp and to_date skip multiple blank spaces at the beginning of the input string and around date and time values unless the FX option is used. For example, to_timestamp_p(' 2000 JUN', 'YYYY MON') and to_timestamp('2000 - JUN', 'YYYY-MON') work, but to_timestamp('2000 JUN', 'FXYYYY MON') returns an error because to_timestamp expects only a single space. FX must be specified as the first item in the template.
- A separator (a space or non-letter/non-digit character) in the template string of to_timestamp and to_date matches any single separator in the input string or is skipped, unless the FX option is used. For example, to_timestamp('2000JUN', 'YYYY//MON') and to_timestamp('2000/JUN', 'YYYY MON') work, but to_timestamp('2000//JUN', 'YYYY/MON') returns an error because the number of separators in the input string exceeds the number of separators in the template.

If FX is specified, a separator in the template string matches exactly one character in the input string. But note that the input string character is not required to be the same as the separator from the template string. For example, to_timestamp('2000/JUN', 'FXYYYY MON') works, but to_timestamp('2000/JUN', 'FXYYYY MON') returns an error because the second space in the template string consumes the letter J from the input string.

- A TZh template pattern can match a signed number. Without the FX option, minus signs may be ambiguous, and could be interpreted as a separator. This ambiguity is resolved as follows: If the number of separators before TZh in the template string is less than the number of separators before the minus sign in the input string, the minus sign is interpreted as part of TZh. Otherwise, the minus sign is considered to be a separator between values. For example, to_timestamp('2000 -10', 'YYYY TZh') matches -10 to TZh, but to_timestamp('2000 -10', 'YYYY TZh') matches 10 to TZh.
- Ordinary text is allowed in to_char templates and will be output literally. You can put a substring in double quotes to force it to be interpreted as literal text even if it contains template patterns. For example, in "Hello Year "YYYY", the YYYY will be replaced by the year data, but the single Y in Year will not be. In to_date, to_number, and to_timestamp, literal text and double-quoted strings result in skipping the number of characters contained in the string; for example "XX" skips two input characters (whether or not they are XX).

Tip

Prior to Postgres Pro 12, it was possible to skip arbitrary text in the input string using non-letter or non-digit characters. For example, `to_timestamp('2000y6m1d', 'yyyy-MM-DD')` used to work. Now you can only use letter characters for this purpose. For example, `to_timestamp('2000y6m1d', 'yyyytMMtDDt')` and `to_timestamp('2000y6m1d', 'yyyy"y"MM"m"DD"d")` skip `y`, `m`, and `d`.

- If you want to have a double quote in the output you must precede it with a backslash, for example `'\"YYYY Month\"'`. Backslashes are not otherwise special outside of double-quoted strings. Within a double-quoted string, a backslash causes the next character to be taken literally, whatever it is (but this has no special effect unless the next character is a double quote or another backslash).
- In `to_timestamp` and `to_date`, if the year format specification is less than four digits, e.g., `YYY`, and the supplied year is less than four digits, the year will be adjusted to be nearest to the year 2020, e.g., `95` becomes 1995.
- In `to_timestamp` and `to_date`, negative years are treated as signifying BC. If you write both a negative year and an explicit BC field, you get AD again. An input of year zero is treated as 1 BC.
- In `to_timestamp` and `to_date`, the `YYYY` conversion has a restriction when processing years with more than 4 digits. You must use some non-digit character or template after `YYYY`, otherwise the year is always interpreted as 4 digits. For example (with the year 20000): `to_date('200001130', 'YYYYMMDD')` will be interpreted as a 4-digit year; instead use a non-digit separator after the year, like `to_date('20000-1130', 'YYYY-MMDD')` or `to_date('20000Nov30', 'YYYYMonDD')`.
- In `to_timestamp` and `to_date`, the `CC` (century) field is accepted but ignored if there is a `YYY`, `YYYY` or `Y`, `YYY` field. If `CC` is used with `YY` or `Y` then the result is computed as that year in the specified century. If the century is specified but the year is not, the first year of the century is assumed.
- In `to_timestamp` and `to_date`, weekday names or numbers (`DAY`, `D`, and related field types) are accepted but are ignored for purposes of computing the result. The same is true for quarter (`Q`) fields.
- In `to_timestamp` and `to_date`, an ISO 8601 week-numbering date (as distinct from a Gregorian date) can be specified in one of two ways:
 - Year, week number, and weekday: for example `to_date('2006-42-4', 'IYYY-IW-ID')` returns the date 2006-10-19. If you omit the weekday it is assumed to be 1 (Monday).
 - Year and day of year: for example `to_date('2006-291', 'IYYY-IDDD')` also returns 2006-10-19.

Attempting to enter a date using a mixture of ISO 8601 week-numbering fields and Gregorian date fields is nonsensical, and will cause an error. In the context of an ISO 8601 week-numbering year, the concept of a “month” or “day of month” has no meaning. In the context of a Gregorian year, the ISO week has no meaning.

Caution

While `to_date` will reject a mixture of Gregorian and ISO week-numbering date fields, `to_char` will not, since output format specifications like `YYYY-MM-DD` (`IYYY-IDDD`) can be useful. But avoid writing something like `IYYY-MM-DD`; that would yield surprising results near the start of the year. (See [Section 9.9.1](#) for more information.)

- In `to_timestamp`, millisecond (`MS`) or microsecond (`US`) fields are used as the seconds digits after the decimal point. For example `to_timestamp('12.3', 'SS.MS')` is not 3 milliseconds, but 300, because the conversion treats it as 12 + 0.3 seconds. So, for the format `SS.MS`, the input values

12.3, 12.30, and 12.300 specify the same number of milliseconds. To get three milliseconds, one must write 12.003, which the conversion treats as $12 + 0.003 = 12.003$ seconds.

Here is a more complex example: `to_timestamp('15:12:02.020.001230', 'HH24:MI:SS.MS.US')` is 15 hours, 12 minutes, and 2 seconds + 20 milliseconds + 1230 microseconds = 2.021230 seconds.

- `to_char(..., 'ID')`'s day of the week numbering matches the `extract(isodow from ...)` function, but `to_char(..., 'D')`'s does not match `extract(dow from ...)`'s day numbering.
- `to_char(interval)` formats HH and HH12 as shown on a 12-hour clock, for example zero hours and 36 hours both output as 12, while HH24 outputs the full hour value, which can exceed 23 in an interval value.

Table 9.29 shows the template patterns available for formatting numeric values.

Table 9.29. Template Patterns for Numeric Formatting

Pattern	Description
9	digit position (can be dropped if insignificant)
0	digit position (will not be dropped, even if insignificant)
. (period)	decimal point
, (comma)	group (thousands) separator
PR	negative value in angle brackets
S	sign anchored to number (uses locale)
L	currency symbol (uses locale)
D	decimal point (uses locale)
G	group separator (uses locale)
MI	minus sign in specified position (if number < 0)
PL	plus sign in specified position (if number > 0)
SG	plus/minus sign in specified position
RN	Roman numeral (input between 1 and 3999)
TH or th	ordinal number suffix
V	shift specified number of digits (see notes)
EEEE	exponent for scientific notation

Usage notes for numeric formatting:

- 0 specifies a digit position that will always be printed, even if it contains a leading/trailing zero. 9 also specifies a digit position, but if it is a leading zero then it will be replaced by a space, while if it is a trailing zero and fill mode is specified then it will be deleted. (For `to_number()`, these two pattern characters are equivalent.)
- If the format provides fewer fractional digits than the number being formatted, `to_char()` will round the number to the specified number of fractional digits.
- The pattern characters S, L, D, and G represent the sign, currency symbol, decimal point, and thousands separator characters defined by the current locale (see [lc_monetary](#) and [lc_numeric](#)). The pattern characters period and comma represent those exact characters, with the meanings of decimal point and thousands separator, regardless of locale.
- If no explicit provision is made for a sign in `to_char()`'s pattern, one column will be reserved for the sign, and it will be anchored to (appear just left of) the number. If S appears just left of some 9's, it will likewise be anchored to the number.

- A sign formatted using `SG`, `PL`, or `MI` is not anchored to the number; for example, `to_char(-12, 'MI9999')` produces `'- 12'` but `to_char(-12, 'S9999')` produces `' -12'`. (The Oracle implementation does not allow the use of `MI` before 9, but rather requires that 9 precede `MI`.)
- `TH` does not convert values less than zero and does not convert fractional numbers.
- `PL`, `SG`, and `TH` are Postgres Pro extensions.
- In `to_number`, if non-data template patterns such as `L` or `TH` are used, the corresponding number of input characters are skipped, whether or not they match the template pattern, unless they are data characters (that is, digits, sign, decimal point, or comma). For example, `TH` would skip two non-data characters.
- `V` with `to_char` multiplies the input values by 10^n , where n is the number of digits following `V`. `V` with `to_number` divides in a similar manner. `to_char` and `to_number` do not support the use of `V` combined with a decimal point (e.g., `99.9V99` is not allowed).
- `EEEE` (scientific notation) cannot be used in combination with any of the other formatting patterns or modifiers other than digit and decimal point patterns, and must be at the end of the format string (e.g., `9.99EEEE` is a valid pattern).

Certain modifiers can be applied to any template pattern to alter its behavior. For example, `FM99.99` is the `99.99` pattern with the `FM` modifier. Table 9.30 shows the modifier patterns for numeric formatting.

Table 9.30. Template Pattern Modifiers for Numeric Formatting

Modifier	Description	Example
FM prefix	fill mode (suppress trailing zeroes and padding blanks)	<code>FM99.99</code>
TH suffix	upper case ordinal number suffix	<code>999TH</code>
th suffix	lower case ordinal number suffix	<code>999th</code>

Table 9.31 shows some examples of the use of the `to_char` function.

Table 9.31. to_char Examples

Expression	Result
<code>to_char(current_timestamp, 'Day, DD HH12:MI:SS')</code>	<code>'Tuesday , 06 05:39:18'</code>
<code>to_char(current_timestamp, 'FMDay, FMDD HH12:MI:SS')</code>	<code>'Tuesday, 6 05:39:18'</code>
<code>to_char(current_timestamp AT TIME ZONE 'UTC', 'YYYY-MM-DD"T"HH24:MI:SS"Z")</code>	<code>'2022-12-06T05:39:18Z', ISO 8601 extended format</code>
<code>to_char(-0.1, '99.99')</code>	<code>' -.10'</code>
<code>to_char(-0.1, 'FM9.99')</code>	<code>'-.1'</code>
<code>to_char(-0.1, 'FM90.99')</code>	<code>'-0.1'</code>
<code>to_char(0.1, '0.9')</code>	<code>' 0.1'</code>
<code>to_char(12, '9990999.9')</code>	<code>' 0012.0'</code>
<code>to_char(12, 'FM9990999.9')</code>	<code>'0012.'</code>
<code>to_char(485, '999')</code>	<code>' 485'</code>
<code>to_char(-485, '999')</code>	<code>'-485'</code>
<code>to_char(485, '9 9 9')</code>	<code>' 4 8 5'</code>
<code>to_char(1485, '9,999')</code>	<code>' 1,485'</code>

Expression	Result
<code>to_char(1485, '9G999')</code>	' 1 485'
<code>to_char(148.5, '999.999')</code>	' 148.500'
<code>to_char(148.5, 'FM999.999')</code>	'148.5'
<code>to_char(148.5, 'FM999.990')</code>	'148.500'
<code>to_char(148.5, '999D999')</code>	' 148,500'
<code>to_char(3148.5, '9G999D999')</code>	' 3 148,500'
<code>to_char(-485, '999S')</code>	'485-'
<code>to_char(-485, '999MI')</code>	'485-'
<code>to_char(485, '999MI')</code>	'485 '
<code>to_char(485, 'FM999MI')</code>	'485'
<code>to_char(485, 'PL999')</code>	'+485'
<code>to_char(485, 'SG999')</code>	'+485'
<code>to_char(-485, 'SG999')</code>	'-485'
<code>to_char(-485, '9SG99')</code>	'4-85'
<code>to_char(-485, '999PR')</code>	'<485>'
<code>to_char(485, 'L999')</code>	'DM 485'
<code>to_char(485, 'RN')</code>	' CDLXXXV'
<code>to_char(485, 'FMRN')</code>	'CDLXXXV'
<code>to_char(5.2, 'FMRN')</code>	'V'
<code>to_char(482, '999th')</code>	' 482nd'
<code>to_char(485, '"Good number:"999')</code>	'Good number: 485'
<code>to_char(485.8, '"Pre:"999" Post:" .999')</code>	'Pre: 485 Post: .800'
<code>to_char(12, '99V999')</code>	' 12000'
<code>to_char(12.4, '99V999')</code>	' 12400'
<code>to_char(12.45, '99V9')</code>	' 125'
<code>to_char(0.0004859, '9.99EEEE')</code>	' 4.86e-04'

9.9. Date/Time Functions and Operators

[Table 9.33](#) shows the available functions for date/time value processing, with details appearing in the following subsections. [Table 9.32](#) illustrates the behaviors of the basic arithmetic operators (+, *, etc.). For formatting functions, refer to [Section 9.8](#). You should be familiar with the background information on date/time data types from [Section 8.5](#).

In addition, the usual comparison operators shown in [Table 9.1](#) are available for the date/time types. Dates and timestamps (with or without time zone) are all comparable, while times (with or without time zone) and intervals can only be compared to other values of the same data type. When comparing a timestamp without time zone to a timestamp with time zone, the former value is assumed to be given in the time zone specified by the [TimeZone](#) configuration parameter, and is rotated to UTC for comparison to the latter value (which is already in UTC internally). Similarly, a date value is assumed to represent midnight in the [TimeZone](#) zone when comparing it to a timestamp.

All the functions and operators described below that take `time` or `timestamp` inputs actually come in two variants: one that takes `time with time zone` or `timestamp with time zone`, and one that takes `time without time zone` or `timestamp without time zone`. For brevity, these variants are not shown

separately. Also, the + and * operators come in commutative pairs (for example both `date + integer` and `integer + date`); we show only one of each such pair.

Table 9.32. Date/Time Operators

Operator Description Example(s)
<code>date + integer → date</code> Add a number of days to a date <code>date '2001-09-28' + 7 → 2001-10-05</code>
<code>date + interval → timestamp</code> Add an interval to a date <code>date '2001-09-28' + interval '1 hour' → 2001-09-28 01:00:00</code>
<code>date + time → timestamp</code> Add a time-of-day to a date <code>date '2001-09-28' + time '03:00' → 2001-09-28 03:00:00</code>
<code>interval + interval → interval</code> Add intervals <code>interval '1 day' + interval '1 hour' → 1 day 01:00:00</code>
<code>timestamp + interval → timestamp</code> Add an interval to a timestamp <code>timestamp '2001-09-28 01:00' + interval '23 hours' → 2001-09-29 00:00:00</code>
<code>time + interval → time</code> Add an interval to a time <code>time '01:00' + interval '3 hours' → 04:00:00</code>
<code>- interval → interval</code> Negate an interval <code>- interval '23 hours' → -23:00:00</code>
<code>date - date → integer</code> Subtract dates, producing the number of days elapsed <code>date '2001-10-01' - date '2001-09-28' → 3</code>
<code>date - integer → date</code> Subtract a number of days from a date <code>date '2001-10-01' - 7 → 2001-09-24</code>
<code>date - interval → timestamp</code> Subtract an interval from a date <code>date '2001-09-28' - interval '1 hour' → 2001-09-27 23:00:00</code>
<code>time - time → interval</code> Subtract times <code>time '05:00' - time '03:00' → 02:00:00</code>
<code>time - interval → time</code> Subtract an interval from a time <code>time '05:00' - interval '2 hours' → 03:00:00</code>
<code>timestamp - interval → timestamp</code> Subtract an interval from a timestamp <code>timestamp '2001-09-28 23:00' - interval '23 hours' → 2001-09-28 00:00:00</code>

Operator	Description	Example(s)
<code>interval - interval</code>	<code>→ interval</code> Subtract intervals	<code>interval '1 day' - interval '1 hour' → 1 day -01:00:00</code>
<code>timestamp - timestamp</code>	<code>→ interval</code> Subtract timestamps (converting 24-hour intervals into days, similarly to <code>justify_hours()</code>)	<code>timestamp '2001-09-29 03:00' - timestamp '2001-07-27 12:00' → 63 days 15:00:00</code>
<code>interval * double precision</code>	<code>→ interval</code> Multiply an interval by a scalar	<code>interval '1 second' * 900 → 00:15:00</code> <code>interval '1 day' * 21 → 21 days</code> <code>interval '1 hour' * 3.5 → 03:30:00</code>
<code>interval / double precision</code>	<code>→ interval</code> Divide an interval by a scalar	<code>interval '1 hour' / 1.5 → 00:40:00</code>

Table 9.33. Date/Time Functions

Function	Description	Example(s)
<code>age (timestamp, timestamp)</code>	<code>→ interval</code> Subtract arguments, producing a “symbolic” result that uses years and months, rather than just days	<code>age(timestamp '2001-04-10', timestamp '1957-06-13') → 43 years 9 mons 27 days</code>
<code>age (timestamp)</code>	<code>→ interval</code> Subtract argument from <code>current_date</code> (at midnight)	<code>age(timestamp '1957-06-13') → 62 years 6 mons 10 days</code>
<code>clock_timestamp ()</code>	<code>→ timestamp with time zone</code> Current date and time (changes during statement execution); see Section 9.9.5	<code>clock_timestamp() → 2019-12-23 14:39:53.662522-05</code>
<code>current_date</code>	<code>→ date</code> Current date; see Section 9.9.5	<code>current_date → 2019-12-23</code>
<code>current_time</code>	<code>→ time with time zone</code> Current time of day; see Section 9.9.5	<code>current_time → 14:39:53.662522-05</code>
<code>current_time (integer)</code>	<code>→ time with time zone</code> Current time of day, with limited precision; see Section 9.9.5	<code>current_time(2) → 14:39:53.66-05</code>
<code>current_timestamp</code>	<code>→ timestamp with time zone</code> Current date and time (start of current transaction); see Section 9.9.5	<code>current_timestamp → 2019-12-23 14:39:53.662522-05</code>
<code>current_timestamp (integer)</code>	<code>→ timestamp with time zone</code> Current date and time (start of current transaction), with limited precision; see Section 9.9.5	

Function	Description	Example(s)
		<code>current_timestamp(0)</code> → 2019-12-23 14:39:53-05
	<code>date_add (timestamp with time zone, interval [, text])</code> → timestamp with time zone Add an interval to a timestamp with time zone, computing times of day and daylight-savings adjustments according to the time zone named by the third argument, or the current TimeZone setting if that is omitted. The form with two arguments is equivalent to the timestamp with time zone + interval operator.	<code>date_add('2021-10-31 00:00:00+02'::timestamp, '1 day'::interval, 'Europe/Warsaw')</code> → 2021-10-31 23:00:00+00
	<code>date_bin (interval, timestamp, timestamp)</code> → timestamp Bin input into specified interval aligned with specified origin; see Section 9.9.3	<code>date_bin('15 minutes', timestamp '2001-02-16 20:38:40', timestamp '2001-02-16 20:05:00')</code> → 2001-02-16 20:35:00
	<code>date_part (text, timestamp)</code> → double precision Get timestamp subfield (equivalent to <code>extract</code>); see Section 9.9.1	<code>date_part('hour', timestamp '2001-02-16 20:38:40')</code> → 20
	<code>date_part (text, interval)</code> → double precision Get interval subfield (equivalent to <code>extract</code>); see Section 9.9.1	<code>date_part('month', interval '2 years 3 months')</code> → 3
	<code>date_subtract (timestamp with time zone, interval [, text])</code> → timestamp with time zone Subtract an interval from a timestamp with time zone, computing times of day and daylight-savings adjustments according to the time zone named by the third argument, or the current TimeZone setting if that is omitted. The form with two arguments is equivalent to the timestamp with time zone - interval operator.	<code>date_subtract('2021-11-01 00:00:00+01'::timestamp, '1 day'::interval, 'Europe/Warsaw')</code> → 2021-10-30 22:00:00+00
	<code>date_trunc (text, timestamp)</code> → timestamp Truncate to specified precision; see Section 9.9.2	<code>date_trunc('hour', timestamp '2001-02-16 20:38:40')</code> → 2001-02-16 20:00:00
	<code>date_trunc (text, timestamp with time zone, text)</code> → timestamp with time zone Truncate to specified precision in the specified time zone; see Section 9.9.2	<code>date_trunc('day', timestamp '2001-02-16 20:38:40+00', 'Australia/Sydney')</code> → 2001-02-16 13:00:00+00
	<code>date_trunc (text, interval)</code> → interval Truncate to specified precision; see Section 9.9.2	<code>date_trunc('hour', interval '2 days 3 hours 40 minutes')</code> → 2 days 03:00:00
	<code>extract (field from timestamp)</code> → numeric Get timestamp subfield; see Section 9.9.1	<code>extract(hour from timestamp '2001-02-16 20:38:40')</code> → 20
	<code>extract (field from interval)</code> → numeric Get interval subfield; see Section 9.9.1	<code>extract(month from interval '2 years 3 months')</code> → 3
	<code>isfinite (date)</code> → boolean Test for finite date (not +/-infinity)	<code>isfinite(date '2001-02-16')</code> → true

Function	Description	Example(s)
<code>isfinite (timestamp)</code>	→ boolean Test for finite timestamp (not +/-infinity)	<code>isfinite(timestamp 'infinity') → false</code>
<code>isfinite (interval)</code>	→ boolean Test for finite interval (not +/-infinity)	<code>isfinite(interval '4 hours') → true</code>
<code>justify_days (interval)</code>	→ interval Adjust interval, converting 30-day time periods to months	<code>justify_days(interval '1 year 65 days') → 1 year 2 mons 5 days</code>
<code>justify_hours (interval)</code>	→ interval Adjust interval, converting 24-hour time periods to days	<code>justify_hours(interval '50 hours 10 minutes') → 2 days 02:10:00</code>
<code>justify_interval (interval)</code>	→ interval Adjust interval using <code>justify_days</code> and <code>justify_hours</code> , with additional sign adjustments	<code>justify_interval(interval '1 mon -1 hour') → 29 days 23:00:00</code>
<code>localtime</code>	→ time Current time of day; see Section 9.9.5	<code>localtime → 14:39:53.662522</code>
<code>localtime (integer)</code>	→ time Current time of day, with limited precision; see Section 9.9.5	<code>localtime(0) → 14:39:53</code>
<code>localtimestamp</code>	→ timestamp Current date and time (start of current transaction); see Section 9.9.5	<code>localtimestamp → 2019-12-23 14:39:53.662522</code>
<code>localtimestamp (integer)</code>	→ timestamp Current date and time (start of current transaction), with limited precision; see Section 9.9.5	<code>localtimestamp(2) → 2019-12-23 14:39:53.66</code>
<code>make_date (year int, month int, day int)</code>	→ date Create date from year, month and day fields (negative years signify BC)	<code>make_date(2013, 7, 15) → 2013-07-15</code>
<code>make_interval ([years int [, months int [, weeks int [, days int [, hours int [, mins int [, secs double precision]]]]]])</code>	→ interval Create interval from years, months, weeks, days, hours, minutes and seconds fields, each of which can default to zero	<code>make_interval(days => 10) → 10 days</code>
<code>make_time (hour int, min int, sec double precision)</code>	→ time Create time from hour, minute and seconds fields	<code>make_time(8, 15, 23.5) → 08:15:23.5</code>
<code>make_timestamp (year int, month int, day int, hour int, min int, sec double precision)</code>	→ timestamp Create timestamp from year, month, day, hour, minute and seconds fields (negative years signify BC)	<code>make_timestamp(2013, 7, 15, 8, 15, 23.5) → 2013-07-15 08:15:23.5</code>

Function	Description	Example(s)
<code>make_timestamptz</code>	<code>(year int, month int, day int, hour int, min int, sec double precision [, timezone text])</code> → timestamp with time zone Create timestamp with time zone from year, month, day, hour, minute and seconds fields (negative years signify BC). If <i>timezone</i> is not specified, the current time zone is used; the examples assume the session time zone is Europe/London	<code>make_timestamptz(2013, 7, 15, 8, 15, 23.5)</code> → 2013-07-15 08:15:23.5+01 <code>make_timestamptz(2013, 7, 15, 8, 15, 23.5, 'America/New_York')</code> → 2013-07-15 13:15:23.5+01
<code>now</code>	<code>()</code> → timestamp with time zone Current date and time (start of current transaction); see Section 9.9.5	<code>now()</code> → 2019-12-23 14:39:53.662522-05
<code>statement_timestamp</code>	<code>()</code> → timestamp with time zone Current date and time (start of current statement); see Section 9.9.5	<code>statement_timestamp()</code> → 2019-12-23 14:39:53.662522-05
<code>timeofday</code>	<code>()</code> → text Current date and time (like <code>clock_timestamp</code> , but as a text string); see Section 9.9.5	<code>timeofday()</code> → Mon Dec 23 14:39:53.662522 2019 EST
<code>transaction_timestamp</code>	<code>()</code> → timestamp with time zone Current date and time (start of current transaction); see Section 9.9.5	<code>transaction_timestamp()</code> → 2019-12-23 14:39:53.662522-05
<code>to_timestamp</code>	<code>(double precision)</code> → timestamp with time zone Convert Unix epoch (seconds since 1970-01-01 00:00:00+00) to timestamp with time zone	<code>to_timestamp(1284352323)</code> → 2010-09-13 04:32:03+00

In addition to these functions, the SQL `OVERLAPS` operator is supported:

```
(start1, end1) OVERLAPS (start2, end2)
(start1, length1) OVERLAPS (start2, length2)
```

This expression yields true when two time periods (defined by their endpoints) overlap, false when they do not overlap. The endpoints can be specified as pairs of dates, times, or time stamps; or as a date, time, or time stamp followed by an interval. When a pair of values is provided, either the start or the end can be written first; `OVERLAPS` automatically takes the earlier value of the pair as the start. Each time period is considered to represent the half-open interval *start* ≤ *time* < *end*, unless *start* and *end* are equal in which case it represents that single time instant. This means for instance that two time periods with only an endpoint in common do not overlap.

```
SELECT (DATE '2001-02-16', DATE '2001-12-21') OVERLAPS
      (DATE '2001-10-30', DATE '2002-10-30');
Result: true
SELECT (DATE '2001-02-16', INTERVAL '100 days') OVERLAPS
      (DATE '2001-10-30', DATE '2002-10-30');
Result: false
SELECT (DATE '2001-10-29', DATE '2001-10-30') OVERLAPS
      (DATE '2001-10-30', DATE '2001-10-31');
Result: false
SELECT (DATE '2001-10-30', DATE '2001-10-30') OVERLAPS
      (DATE '2001-10-30', DATE '2001-10-31');
Result: true
```

When adding an interval value to (or subtracting an interval value from) a timestamp or timestamp with time zone value, the months, days, and microseconds fields of the interval value are handled in turn. First, a nonzero months field advances or decrements the date of the timestamp by the indicated number of months, keeping the day of month the same unless it would be past the end of the new month, in which case the last day of that month is used. (For example, March 31 plus 1 month becomes April 30, but March 31 plus 2 months becomes May 31.) Then the days field advances or decrements the date of the timestamp by the indicated number of days. In both these steps the local time of day is kept the same. Finally, if there is a nonzero microseconds field, it is added or subtracted literally. When doing arithmetic on a timestamp with time zone value in a time zone that recognizes DST, this means that adding or subtracting (say) interval '1 day' does not necessarily have the same result as adding or subtracting interval '24 hours'. For example, with the session time zone set to America/Denver:

```
SELECT timestamp with time zone '2005-04-02 12:00:00-07' + interval '1 day';
Result: 2005-04-03 12:00:00-06
SELECT timestamp with time zone '2005-04-02 12:00:00-07' + interval '24 hours';
Result: 2005-04-03 13:00:00-06
```

This happens because an hour was skipped due to a change in daylight saving time at 2005-04-03 02:00:00 in time zone America/Denver.

Note there can be ambiguity in the months field returned by age because different months have different numbers of days. Postgres Pro's approach uses the month from the earlier of the two dates when calculating partial months. For example, age('2004-06-01', '2004-04-30') uses April to yield 1 mon 1 day, while using May would yield 1 mon 2 days because May has 31 days, while April has only 30.

Subtraction of dates and timestamps can also be complex. One conceptually simple way to perform subtraction is to convert each value to a number of seconds using EXTRACT(EPOCH FROM ...), then subtract the results; this produces the number of seconds between the two values. This will adjust for the number of days in each month, timezone changes, and daylight saving time adjustments. Subtraction of date or timestamp values with the "-" operator returns the number of days (24-hours) and hours/minutes/seconds between the values, making the same adjustments. The age function returns years, months, days, and hours/minutes/seconds, performing field-by-field subtraction and then adjusting for negative field values. The following queries illustrate the differences in these approaches. The sample results were produced with timezone = 'US/Eastern'; there is a daylight saving time change between the two dates used:

```
SELECT EXTRACT(EPOCH FROM timestampz '2013-07-01 12:00:00') -
       EXTRACT(EPOCH FROM timestampz '2013-03-01 12:00:00');
Result: 10537200.000000
SELECT (EXTRACT(EPOCH FROM timestampz '2013-07-01 12:00:00') -
       EXTRACT(EPOCH FROM timestampz '2013-03-01 12:00:00'))
       / 60 / 60 / 24;
Result: 121.95833333333333
SELECT timestampz '2013-07-01 12:00:00' - timestampz '2013-03-01 12:00:00';
Result: 121 days 23:00:00
SELECT age(timestampz '2013-07-01 12:00:00', timestampz '2013-03-01 12:00:00');
Result: 4 mons
```

9.9.1. EXTRACT, date_part

EXTRACT(*field* FROM *source*)

The extract function retrieves subfields such as year or hour from date/time values. *source* must be a value expression of type timestamp, date, time, or interval. (Timestamps and times can be with or without time zone.) *field* is an identifier or string that selects what field to extract from the source value. Not all fields are valid for every input data type; for example, fields smaller than a day cannot be extracted from a date, while fields of a day or more cannot be extracted from a time. The extract function returns values of type numeric.

The following are valid field names:

century

The century; for interval values, the year field divided by 100

```
SELECT EXTRACT(CENTURY FROM TIMESTAMP '2000-12-16 12:21:13');
Result: 20
SELECT EXTRACT(CENTURY FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 21
SELECT EXTRACT(CENTURY FROM DATE '0001-01-01 AD');
Result: 1
SELECT EXTRACT(CENTURY FROM DATE '0001-12-31 BC');
Result: -1
SELECT EXTRACT(CENTURY FROM INTERVAL '2001 years');
Result: 20
```

day

The day of the month (1-31); for interval values, the number of days

```
SELECT EXTRACT(DAY FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 16
SELECT EXTRACT(DAY FROM INTERVAL '40 days 1 minute');
Result: 40
```

decade

The year field divided by 10

```
SELECT EXTRACT(DECADE FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 200
```

dow

The day of the week as Sunday (0) to Saturday (6)

```
SELECT EXTRACT(DOW FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 5
```

Note that `extract`'s day of the week numbering differs from that of the `to_char(..., 'D')` function.

doy

The day of the year (1-365/366)

```
SELECT EXTRACT(DOY FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 47
```

epoch

For timestamp with time zone values, the number of seconds since 1970-01-01 00:00:00 UTC (negative for timestamps before that); for date and timestamp values, the nominal number of seconds since 1970-01-01 00:00:00, without regard to timezone or daylight-savings rules; for interval values, the total number of seconds in the interval

```
SELECT EXTRACT(EPOCH FROM TIMESTAMP WITH TIME ZONE '2001-02-16 20:38:40.12-08');
Result: 982384720.120000
SELECT EXTRACT(EPOCH FROM TIMESTAMP '2001-02-16 20:38:40.12');
Result: 982355920.120000
SELECT EXTRACT(EPOCH FROM INTERVAL '5 days 3 hours');
Result: 442800.000000
```

You can convert an epoch value back to a timestamp with time zone with `to_timestamp`:

```
SELECT to_timestamp(982384720.12);
Result: 2001-02-17 04:38:40.12+00
```


Beware that applying `to_timestamp` to an epoch extracted from a `date` or `timestamp` value could produce a misleading result: the result will effectively assume that the original value had been given in UTC, which might not be the case.

hour

The hour field (0-23 in timestamps, unrestricted in intervals)

```
SELECT EXTRACT(HOUR FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 20
```

isodow

The day of the week as Monday (1) to Sunday (7)

```
SELECT EXTRACT(ISODOW FROM TIMESTAMP '2001-02-18 20:38:40');  
Result: 7
```

This is identical to `dow` except for Sunday. This matches the ISO 8601 day of the week numbering.

isoyear

The ISO 8601 week-numbering year that the date falls in

```
SELECT EXTRACT(ISOYEAR FROM DATE '2006-01-01');  
Result: 2005  
SELECT EXTRACT(ISOYEAR FROM DATE '2006-01-02');  
Result: 2006
```

Each ISO 8601 week-numbering year begins with the Monday of the week containing the 4th of January, so in early January or late December the ISO year may be different from the Gregorian year. See the `week` field for more information.

julian

The *Julian Date* corresponding to the date or timestamp. Timestamps that are not local midnight result in a fractional value. See [Section B.7](#) for more information.

```
SELECT EXTRACT(JULIAN FROM DATE '2006-01-01');  
Result: 2453737  
SELECT EXTRACT(JULIAN FROM TIMESTAMP '2006-01-01 12:00');  
Result: 2453737.500000000000000000000000
```

microseconds

The seconds field, including fractional parts, multiplied by 1 000 000; note that this includes full seconds

```
SELECT EXTRACT(MICROSECONDS FROM TIME '17:12:28.5');  
Result: 28500000
```

millennium

The millennium; for `interval` values, the year field divided by 1000

```
SELECT EXTRACT(MILLENNIUM FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 3  
SELECT EXTRACT(MILLENNIUM FROM INTERVAL '2001 years');  
Result: 2
```

Years in the 1900s are in the second millennium. The third millennium started January 1, 2001.

milliseconds

The seconds field, including fractional parts, multiplied by 1000. Note that this includes full seconds.

```
SELECT EXTRACT(MILLISECONDS FROM TIME '17:12:28.5');  
Result: 28500.000
```

minute

The minutes field (0-59)

```
SELECT EXTRACT(MINUTE FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 38
```

month

The number of the month within the year (1-12); for interval values, the number of months modulo 12 (0-11)

```
SELECT EXTRACT(MONTH FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 2  
SELECT EXTRACT(MONTH FROM INTERVAL '2 years 3 months');  
Result: 3  
SELECT EXTRACT(MONTH FROM INTERVAL '2 years 13 months');  
Result: 1
```

quarter

The quarter of the year (1-4) that the date is in

```
SELECT EXTRACT(QUARTER FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 1
```

second

The seconds field, including any fractional seconds

```
SELECT EXTRACT(SECOND FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 40.000000  
SELECT EXTRACT(SECOND FROM TIME '17:12:28.5');  
Result: 28.500000
```

timezone

The time zone offset from UTC, measured in seconds. Positive values correspond to time zones east of UTC, negative values to zones west of UTC. (Technically, Postgres Pro does not use UTC because leap seconds are not handled.)

timezone_hour

The hour component of the time zone offset

timezone_minute

The minute component of the time zone offset

week

The number of the ISO 8601 week-numbering week of the year. By definition, ISO weeks start on Mondays and the first week of a year contains January 4 of that year. In other words, the first Thursday of a year is in week 1 of that year.

In the ISO week-numbering system, it is possible for early-January dates to be part of the 52nd or 53rd week of the previous year, and for late-December dates to be part of the first week of the next year. For example, 2005-01-01 is part of the 53rd week of year 2004, and 2006-01-01 is part of the 52nd week of year 2005, while 2012-12-31 is part of the first week of 2013. It's recommended to use the `isoyear` field together with `week` to get consistent results.

```
SELECT EXTRACT(WEEK FROM TIMESTAMP '2001-02-16 20:38:40');  
Result: 7
```

year

The year field. Keep in mind there is no 0 AD, so subtracting BC years from AD years should be done with care.

```
SELECT EXTRACT(YEAR FROM TIMESTAMP '2001-02-16 20:38:40');
Result: 2001
```

When processing an interval value, the `extract` function produces field values that match the interpretation used by the interval output function. This can produce surprising results if one starts with a non-normalized interval representation, for example:

```
SELECT INTERVAL '80 minutes';
Result: 01:20:00
SELECT EXTRACT(MINUTES FROM INTERVAL '80 minutes');
Result: 20
```

Note

When the input value is +/-Infinity, `extract` returns +/-Infinity for monotonically-increasing fields (epoch, julian, year, isoyear, decade, century, and millennium for timestamp inputs; epoch, hour, day, year, decade, century, and millennium for interval inputs). For other fields, NULL is returned. PostgreSQL versions before 9.6 returned zero for all cases of infinite input.

The `extract` function is primarily intended for computational processing. For formatting date/time values for display, see [Section 9.8](#).

The `date_part` function is modeled on the traditional Ingres equivalent to the SQL-standard function `extract`:

```
date_part('field', source)
```

Note that here the *field* parameter needs to be a string value, not a name. The valid field names for `date_part` are the same as for `extract`. For historical reasons, the `date_part` function returns values of type double precision. This can result in a loss of precision in certain uses. Using `extract` is recommended instead.

```
SELECT date_part('day', TIMESTAMP '2001-02-16 20:38:40');
Result: 16
SELECT date_part('hour', INTERVAL '4 hours 3 minutes');
Result: 4
```

9.9.2. date_trunc

The function `date_trunc` is conceptually similar to the `trunc` function for numbers.

```
date_trunc(field, source [, time_zone ])
```

source is a value expression of type timestamp, timestamp with time zone, or interval. (Values of type date and time are cast automatically to timestamp or interval, respectively.) *field* selects to which precision to truncate the input value. The return value is likewise of type timestamp, timestamp with time zone, or interval, and it has all fields that are less significant than the selected one set to zero (or one, for day and month).

Valid values for *field* are:

```
microseconds
milliseconds
second
minute
hour
```

day
week
month
quarter
year
decade
century
millennium

When the input value is of type `timestamp with time zone`, the truncation is performed with respect to a particular time zone; for example, truncation to `day` produces a value that is midnight in that zone. By default, truncation is done with respect to the current [TimeZone](#) setting, but the optional `time_zone` argument can be provided to specify a different time zone. The time zone name can be specified in any of the ways described in [Section 8.5.3](#).

A time zone cannot be specified when processing `timestamp without time zone` or `interval` inputs. These are always taken at face value.

Examples (assuming the local time zone is `America/New_York`):

```
SELECT date_trunc('hour', TIMESTAMP '2001-02-16 20:38:40');
Result: 2001-02-16 20:00:00
SELECT date_trunc('year', TIMESTAMP '2001-02-16 20:38:40');
Result: 2001-01-01 00:00:00
SELECT date_trunc('day', TIMESTAMP WITH TIME ZONE '2001-02-16 20:38:40+00');
Result: 2001-02-16 00:00:00-05
SELECT date_trunc('day', TIMESTAMP WITH TIME ZONE '2001-02-16 20:38:40+00', 'Australia/Sydney');
Result: 2001-02-16 08:00:00-05
SELECT date_trunc('hour', INTERVAL '3 days 02:47:33');
Result: 3 days 02:00:00
```

9.9.3. date_bin

The function `date_bin` “bins” the input timestamp into the specified interval (the *stride*) aligned with a specified origin.

`date_bin(stride, source, origin)`

`source` is a value expression of type `timestamp` or `timestamp with time zone`. (Values of type `date` are cast automatically to `timestamp`.) `stride` is a value expression of type `interval`. The return value is likewise of type `timestamp` or `timestamp with time zone`, and it marks the beginning of the bin into which the `source` is placed.

Examples:

```
SELECT date_bin('15 minutes', TIMESTAMP '2020-02-11 15:44:17', TIMESTAMP '2001-01-01');
Result: 2020-02-11 15:30:00
SELECT date_bin('15 minutes', TIMESTAMP '2020-02-11 15:44:17', TIMESTAMP '2001-01-01
00:02:30');
Result: 2020-02-11 15:32:30
```

In the case of full units (1 minute, 1 hour, etc.), it gives the same result as the analogous `date_trunc` call, but the difference is that `date_bin` can truncate to an arbitrary interval.

The `stride` interval must be greater than zero and cannot contain units of month or larger.

9.9.4. AT TIME ZONE and AT LOCAL

The `AT TIME ZONE` operator converts time stamp *without* time zone to/from time stamp *with* time zone, and time with time zone values to different time zones. [Table 9.34](#) shows its variants.

Table 9.34. AT TIME ZONE and AT LOCAL Variants

Operator	Description	Example(s)
timestamp without time zone	AT TIME ZONE <i>zone</i> → timestamp with time zone	Converts given time stamp <i>without</i> time zone to time stamp <i>with</i> time zone, assuming the given value is in the named time zone. timestamp '2001-02-16 20:38:40' at time zone 'America/Denver' → 2001-02-17 03:38:40+00
timestamp without time zone	AT LOCAL → timestamp with time zone	Converts given time stamp <i>without</i> time zone to time stamp <i>with</i> the session's TimeZone value as time zone. timestamp '2001-02-16 20:38:40' at local → 2001-02-17 03:38:40+00
timestamp with time zone	AT TIME ZONE <i>zone</i> → timestamp without time zone	Converts given time stamp <i>with</i> time zone to time stamp <i>without</i> time zone, as the time would appear in that zone. timestamp with time zone '2001-02-16 20:38:40-05' at time zone 'America/Denver' → 2001-02-16 18:38:40
timestamp with time zone	AT LOCAL → timestamp without time zone	Converts given time stamp <i>with</i> time zone to time stamp <i>without</i> time zone, as the time would appear with the session's TimeZone value as time zone. timestamp with time zone '2001-02-16 20:38:40-05' at local → 2001-02-16 18:38:40
time with time zone	AT TIME ZONE <i>zone</i> → time with time zone	Converts given time <i>with</i> time zone to a new time zone. Since no date is supplied, this uses the currently active UTC offset for the named destination zone. time with time zone '05:34:17-05' at time zone 'UTC' → 10:34:17+00
time with time zone	AT LOCAL → time with time zone	Converts given time <i>with</i> time zone to a new time zone. Since no date is supplied, this uses the currently active UTC offset for the session's TimeZone value. Assuming the session's TimeZone is set to UTC: time with time zone '05:34:17-05' at local → 10:34:17+00

In these expressions, the desired time zone *zone* can be specified either as a text value (e.g., 'America/Los_Angeles') or as an interval (e.g., INTERVAL '-08:00'). In the text case, a time zone name can be specified in any of the ways described in [Section 8.5.3](#). The interval case is only useful for zones that have fixed offsets from UTC, so it is not very common in practice.

The syntax AT LOCAL may be used as shorthand for AT TIME ZONE *local*, where *local* is the session's TimeZone value.

Examples (assuming the current TimeZone setting is America/Los_Angeles):

```
SELECT TIMESTAMP '2001-02-16 20:38:40' AT TIME ZONE 'America/Denver';
Result: 2001-02-16 19:38:40-08
SELECT TIMESTAMP WITH TIME ZONE '2001-02-16 20:38:40-05' AT TIME ZONE 'America/Denver';
Result: 2001-02-16 18:38:40
SELECT TIMESTAMP '2001-02-16 20:38:40' AT TIME ZONE 'Asia/Tokyo' AT TIME ZONE 'America/Chicago';
Result: 2001-02-16 05:38:40
SELECT TIMESTAMP WITH TIME ZONE '2001-02-16 20:38:40-05' AT LOCAL;
Result: 2001-02-16 17:38:40
SELECT TIME WITH TIME ZONE '20:38:40-05' AT LOCAL;
```

Result: 17:38:40

The first example adds a time zone to a value that lacks it, and displays the value using the current `TimeZone` setting. The second example shifts the time stamp with time zone value to the specified time zone, and returns the value without a time zone. This allows storage and display of values different from the current `TimeZone` setting. The third example converts Tokyo time to Chicago time. The fourth example shifts the time stamp with time zone value to the time zone currently specified by the `TimeZone` setting and returns the value without a time zone.

The fifth example is a cautionary tale. Due to the fact that there is no date associated with the input value, the conversion is made using the current date of the session. Therefore, this static example may show a wrong result depending on the time of the year it is viewed because 'America/Los_Angeles' observes Daylight Savings Time.

The function `timezone(zone, timestamp)` is equivalent to the SQL-conforming construct `timestamp AT TIME ZONE zone`.

The function `timezone(zone, time)` is equivalent to the SQL-conforming construct `time AT TIME ZONE zone`.

The function `timezone(timestamp)` is equivalent to the SQL-conforming construct `timestamp AT LOCAL`.

The function `timezone(time)` is equivalent to the SQL-conforming construct `time AT LOCAL`.

9.9.5. Current Date/Time

Postgres Pro provides a number of functions that return values related to the current date and time. These SQL-standard functions all return values based on the start time of the current transaction:

```
CURRENT_DATE
CURRENT_TIME
CURRENT_TIMESTAMP
CURRENT_TIME(precision)
CURRENT_TIMESTAMP(precision)
LOCALTIME
LOCALTIMESTAMP
LOCALTIME(precision)
LOCALTIMESTAMP(precision)
```

`CURRENT_TIME` and `CURRENT_TIMESTAMP` deliver values with time zone; `LOCALTIME` and `LOCALTIMESTAMP` deliver values without time zone.

`CURRENT_TIME`, `CURRENT_TIMESTAMP`, `LOCALTIME`, and `LOCALTIMESTAMP` can optionally take a precision parameter, which causes the result to be rounded to that many fractional digits in the seconds field. Without a precision parameter, the result is given to the full available precision.

Some examples:

```
SELECT CURRENT_TIME;
Result: 14:39:53.662522-05
SELECT CURRENT_DATE;
Result: 2019-12-23
SELECT CURRENT_TIMESTAMP;
Result: 2019-12-23 14:39:53.662522-05
SELECT CURRENT_TIMESTAMP(2);
Result: 2019-12-23 14:39:53.66-05
SELECT LOCALTIMESTAMP;
Result: 2019-12-23 14:39:53.662522
```

Since these functions return the start time of the current transaction, their values do not change during the transaction. This is considered a feature: the intent is to allow a single transaction to have a consistent

notion of the “current” time, so that multiple modifications within the same transaction bear the same time stamp.

Note

Other database systems might advance these values more frequently.

Postgres Pro also provides functions that return the start time of the current statement, as well as the actual current time at the instant the function is called. The complete list of non-SQL-standard time functions is:

```
transaction_timestamp()
statement_timestamp()
clock_timestamp()
timeofday()
now()
```

`transaction_timestamp()` is equivalent to `CURRENT_TIMESTAMP`, but is named to clearly reflect what it returns. `statement_timestamp()` returns the start time of the current statement (more specifically, the time of receipt of the latest command message from the client). `statement_timestamp()` and `transaction_timestamp()` return the same value during the first command of a transaction, but might differ during subsequent commands. `clock_timestamp()` returns the actual current time, and therefore its value changes even within a single SQL command. `timeofday()` is a historical Postgres Pro function. Like `clock_timestamp()`, it returns the actual current time, but as a formatted text string rather than a timestamp with time zone value. `now()` is a traditional Postgres Pro equivalent to `transaction_timestamp()`.

All the date/time data types also accept the special literal value `now` to specify the current date and time (again, interpreted as the transaction start time). Thus, the following three all return the same result:

```
SELECT CURRENT_TIMESTAMP;
SELECT now();
SELECT TIMESTAMP 'now'; -- but see tip below
```

Tip

Do not use the third form when specifying a value to be evaluated later, for example in a `DEFAULT` clause for a table column. The system will convert `now` to a `timestamp` as soon as the constant is parsed, so that when the default value is needed, the time of the table creation would be used! The first two forms will not be evaluated until the default value is used, because they are function calls. Thus they will give the desired behavior of defaulting to the time of row insertion. (See also [Section 8.5.1.4.](#))

9.9.6. Delaying Execution

The following functions are available to delay execution of the server process:

```
pg_sleep ( double precision )
pg_sleep_for ( interval )
pg_sleep_until ( timestamp with time zone )
```

`pg_sleep` makes the current session's process sleep until the given number of seconds have elapsed. Fractional-second delays can be specified. `pg_sleep_for` is a convenience function to allow the sleep time to be specified as an `interval`. `pg_sleep_until` is a convenience function for when a specific wake-up time is desired. For example:

```
SELECT pg_sleep(1.5);
SELECT pg_sleep_for('5 minutes');
SELECT pg_sleep_until('tomorrow 03:00');
```

Note

The effective resolution of the sleep interval is platform-specific; 0.01 seconds is a common value. The sleep delay will be at least as long as specified. It might be longer depending on factors such as server load. In particular, `pg_sleep_until` is not guaranteed to wake up exactly at the specified time, but it will not wake up any earlier.

Warning

Make sure that your session does not hold more locks than necessary when calling `pg_sleep` or its variants. Otherwise other sessions might have to wait for your sleeping process, slowing down the entire system.

9.10. Enum Support Functions

For enum types (described in [Section 8.7](#)), there are several functions that allow cleaner programming without hard-coding particular values of an enum type. These are listed in [Table 9.35](#). The examples assume an enum type created as:

```
CREATE TYPE rainbow AS ENUM ('red', 'orange', 'yellow', 'green', 'blue', 'purple');
```

Table 9.35. Enum Support Functions

Function	Description	Example(s)
<code>enum_first (anyenum) → anyenum</code>	Returns the first value of the input enum type.	<code>enum_first(null::rainbow) → red</code>
<code>enum_last (anyenum) → anyenum</code>	Returns the last value of the input enum type.	<code>enum_last(null::rainbow) → purple</code>
<code>enum_range (anyenum) → anyarray</code>	Returns all values of the input enum type in an ordered array.	<code>enum_range(null::rainbow) → {red,orange,yellow,green,blue,purple}</code>
<code>enum_range (anyenum, anyenum) → anyarray</code>	Returns the range between the two given enum values, as an ordered array. The values must be from the same enum type. If the first parameter is null, the result will start with the first value of the enum type. If the second parameter is null, the result will end with the last value of the enum type.	<code>enum_range('orange'::rainbow, 'green'::rainbow) → {orange,yellow,green}</code> <code>enum_range(NULL, 'green'::rainbow) → {red,orange,yellow,green}</code> <code>enum_range('orange'::rainbow, NULL) → {orange,yellow,green,blue,purple}</code>

Notice that except for the two-argument form of `enum_range`, these functions disregard the specific value passed to them; they care only about its declared data type. Either null or a specific value of the type can be passed, with the same result. It is more common to apply these functions to a table column or function argument than to a hardwired type name as used in the examples.

9.11. Geometric Functions and Operators

The geometric types `point`, `box`, `lseg`, `line`, `path`, `polygon`, and `circle` have a large set of native support functions and operators, shown in [Table 9.36](#), [Table 9.37](#), and [Table 9.38](#).

Table 9.36. Geometric Operators

Operator	Description	Example(s)
<code>geometric_type + point → geometric_type</code>	Adds the coordinates of the second <code>point</code> to those of each point of the first argument, thus performing translation. Available for <code>point</code> , <code>box</code> , <code>path</code> , <code>circle</code> .	<code>box '(1,1), (0,0)' + point '(2,0)'</code> → <code>(3,1), (2,0)</code>
<code>path + path → path</code>	Concatenates two open paths (returns NULL if either path is closed).	<code>path '[(0,0), (1,1)]' + path '[(2,2), (3,3), (4,4)]'</code> → <code>[(0,0), (1,1), (2,2), (3,3), (4,4)]</code>
<code>geometric_type - point → geometric_type</code>	Subtracts the coordinates of the second <code>point</code> from those of each point of the first argument, thus performing translation. Available for <code>point</code> , <code>box</code> , <code>path</code> , <code>circle</code> .	<code>box '(1,1), (0,0)' - point '(2,0)'</code> → <code>(-1,1), (-2,0)</code>
<code>geometric_type * point → geometric_type</code>	Multiplies each point of the first argument by the second <code>point</code> (treating a point as being a complex number represented by real and imaginary parts, and performing standard complex multiplication). If one interprets the second <code>point</code> as a vector, this is equivalent to scaling the object's size and distance from the origin by the length of the vector, and rotating it counter-clockwise around the origin by the vector's angle from the <code>x</code> axis. Available for <code>point</code> , <code>box</code> , ^a <code>path</code> , <code>circle</code> .	<code>path '((0,0), (1,0), (1,1))' * point '(3.0,0)'</code> → <code>((0,0), (3,0), (3,3))</code> <code>path '((0,0), (1,0), (1,1))' * point (cosd(45), sind(45))</code> → <code>((0,0), (0.7071067811865475, 0.7071067811865475), (0, 1.414213562373095))</code>
<code>geometric_type / point → geometric_type</code>	Divides each point of the first argument by the second <code>point</code> (treating a point as being a complex number represented by real and imaginary parts, and performing standard complex division). If one interprets the second <code>point</code> as a vector, this is equivalent to scaling the object's size and distance from the origin down by the length of the vector, and rotating it clockwise around the origin by the vector's angle from the <code>x</code> axis. Available for <code>point</code> , <code>box</code> , ^a <code>path</code> , <code>circle</code> .	<code>path '((0,0), (1,0), (1,1))' / point '(2.0,0)'</code> → <code>((0,0), (0.5,0), (0.5,0.5))</code> <code>path '((0,0), (1,0), (1,1))' / point (cosd(45), sind(45))</code> → <code>((0,0), (0.7071067811865476, -0.7071067811865476), (1.4142135623730951, 0))</code>
<code>@-@ geometric_type → double precision</code>	Computes the total length. Available for <code>lseg</code> , <code>path</code> .	<code>@-@ path '[(0,0), (1,0), (1,1)]'</code> → <code>2</code>
<code>@@ geometric_type → point</code>	Computes the center point. Available for <code>box</code> , <code>lseg</code> , <code>polygon</code> , <code>circle</code> .	<code>@@ box '(2,2), (0,0)'</code> → <code>(1,1)</code>
<code># geometric_type → integer</code>	Returns the number of points. Available for <code>path</code> , <code>polygon</code> .	<code># path '((1,0), (0,1), (-1,0))'</code> → <code>3</code>
<code>geometric_type # geometric_type → point</code>		

Operator	Description	Example(s)
	Computes the point of intersection, or NULL if there is none. Available for lseg, line.	lseg '(0,0),(1,1)' # lseg '(1,0),(0,1)' → (0.5,0.5)
box # box → box	Computes the intersection of two boxes, or NULL if there is none.	box '(2,2),(-1,-1)' # box '(1,1),(-2,-2)' → (1,1),(-1,-1)
geometric_type ## geometric_type → point	Computes the closest point to the first object on the second object. Available for these pairs of types: (point, box), (point, lseg), (point, line), (lseg, box), (lseg, lseg), (line, lseg).	point '(0,0)' ## lseg '(2,0),(0,2)' → (1,1)
geometric_type <-> geometric_type → double precision	Computes the distance between the objects. Available for all seven geometric types, for all combinations of point with another geometric type, and for these additional pairs of types: (box, lseg), (lseg, line), (polygon, circle) (and the commutator cases).	circle '<(0,0),1>' <-> circle '<(5,0),1>' → 3
geometric_type @> geometric_type → boolean	Does first object contain second? Available for these pairs of types: (box, point), (box, box), (path, point), (polygon, point), (polygon, polygon), (circle, point), (circle, circle).	circle '<(0,0),2>' @> point '(1,1)' → t
geometric_type <@ geometric_type → boolean	Is first object contained in or on second? Available for these pairs of types: (point, box), (point, lseg), (point, line), (point, path), (point, polygon), (point, circle), (box, box), (lseg, box), (lseg, line), (polygon, polygon), (circle, circle).	point '(1,1)' <@ circle '<(0,0),2>' → t
geometric_type && geometric_type → boolean	Do these objects overlap? (One point in common makes this true.) Available for box, polygon, circle.	box '(1,1),(0,0)' && box '(2,2),(0,0)' → t
geometric_type << geometric_type → boolean	Is first object strictly left of second? Available for point, box, polygon, circle.	circle '<(0,0),1>' << circle '<(5,0),1>' → t
geometric_type >> geometric_type → boolean	Is first object strictly right of second? Available for point, box, polygon, circle.	circle '<(5,0),1>' >> circle '<(0,0),1>' → t
geometric_type &< geometric_type → boolean	Does first object not extend to the right of second? Available for box, polygon, circle.	box '(1,1),(0,0)' &< box '(2,2),(0,0)' → t
geometric_type &> geometric_type → boolean	Does first object not extend to the left of second? Available for box, polygon, circle.	box '(3,3),(0,0)' &> box '(2,2),(0,0)' → t
geometric_type << geometric_type → boolean	Is first object strictly below second? Available for point, box, polygon, circle.	box '(3,3),(0,0)' << box '(5,5),(3,4)' → t
geometric_type >> geometric_type → boolean		

Operator	Description	Example(s)
	Is first object strictly above second? Available for point, box, polygon, circle.	<code>box '(5,5), (3,4)' >> box '(3,3), (0,0)' → t</code>
	<code>geometric_type &< geometric_type → boolean</code> Does first object not extend above second? Available for box, polygon, circle.	<code>box '(1,1), (0,0)' &< box '(2,2), (0,0)' → t</code>
	<code>geometric_type &> geometric_type → boolean</code> Does first object not extend below second? Available for box, polygon, circle.	<code>box '(3,3), (0,0)' &> box '(2,2), (0,0)' → t</code>
	<code>box <^ box → boolean</code> Is first object below second (allows edges to touch)?	<code>box '((1,1), (0,0))' <^ box '((2,2), (1,1))' → t</code>
	<code>box >^ box → boolean</code> Is first object above second (allows edges to touch)?	<code>box '((2,2), (1,1))' >^ box '((1,1), (0,0))' → t</code>
	<code>geometric_type ?# geometric_type → boolean</code> Do these objects intersect? Available for these pairs of types: (box, box), (lseg, box), (lseg, lseg), (lseg, line), (line, box), (line, line), (path, path).	<code>lseg '((-1,0), (1,0))' ?# box '(2,2), (-2,-2)' → t</code>
	<code>?- line → boolean</code> <code>?- lseg → boolean</code> Is line horizontal?	<code>?- lseg '((-1,0), (1,0))' → t</code>
	<code>point ?- point → boolean</code> Are points horizontally aligned (that is, have same y coordinate)?	<code>point '(1,0)' ?- point '(0,0)' → t</code>
	<code>? line → boolean</code> <code>? lseg → boolean</code> Is line vertical?	<code>? lseg '((-1,0), (1,0))' → f</code>
	<code>point ? point → boolean</code> Are points vertically aligned (that is, have same x coordinate)?	<code>point '(0,1)' ? point '(0,0)' → t</code>
	<code>line ?- line → boolean</code> <code>lseg ?- lseg → boolean</code> Are lines perpendicular?	<code>lseg '[(0,0), (0,1)]' ?- lseg '[(0,0), (1,0)]' → t</code>
	<code>line ? line → boolean</code> <code>lseg ? lseg → boolean</code> Are lines parallel?	<code>lseg '((-1,0), (1,0))' ? lseg '((-1,2), (1,2))' → t</code>
	<code>geometric_type ~= geometric_type → boolean</code> Are these objects the same? Available for point, box, polygon, circle.	

Operator	Description	Example(s)
		<code>polygon '((0,0),(1,1))' ~= polygon '((1,1),(0,0))'</code> → t

^a“Rotating” a box with these operators only moves its corner points: the box is still considered to have sides parallel to the axes. Hence the box's size is not preserved, as a true rotation would do.

Caution

Note that the “same as” operator, `~=`, represents the usual notion of equality for the `point`, `box`, `polygon`, and `circle` types. Some of the geometric types also have an `=` operator, but `=` compares for equal *areas* only. The other scalar comparison operators (`<=` and so on), where available for these types, likewise compare areas.

Note

Before Postgres Pro 14, the point is strictly below/above comparison operators `point <<| point` and `point |>> point` were respectively called `<^` and `>^`. These names are still available, but are deprecated and will eventually be removed.

Table 9.37. Geometric Functions

Function	Description	Example(s)
<code>area (<i>geometric_type</i>)</code> → double precision	Computes area. Available for <code>box</code> , <code>path</code> , <code>circle</code> . A <code>path</code> input must be closed, else NULL is returned. Also, if the <code>path</code> is self-intersecting, the result may be meaningless.	<code>area(box '(2,2),(0,0)')</code> → 4
<code>center (<i>geometric_type</i>)</code> → point	Computes center point. Available for <code>box</code> , <code>circle</code> .	<code>center(box '(1,2),(0,0)')</code> → (0.5,1)
<code>diagonal (box)</code> → lseg	Extracts <code>box</code> 's diagonal as a line segment (same as <code>lseg(box)</code>).	<code>diagonal(box '(1,2),(0,0)')</code> → [(1,2),(0,0)]
<code>diameter (circle)</code> → double precision	Computes diameter of circle.	<code>diameter(circle '<(0,0),2>')</code> → 4
<code>height (box)</code> → double precision	Computes vertical size of <code>box</code> .	<code>height(box '(1,2),(0,0)')</code> → 2
<code>isclosed (path)</code> → boolean	Is path closed?	<code>isclosed(path '((0,0),(1,1),(2,0))')</code> → t
<code>isopen (path)</code> → boolean	Is path open?	<code>isopen(path '[(0,0),(1,1),(2,0)]')</code> → t
<code>length (<i>geometric_type</i>)</code> → double precision	Computes the total length. Available for <code>lseg</code> , <code>path</code> .	

Function	Description	Example(s)
<code>length (path)</code>	<code>length (path ' ((-1,0) , (1,0)) ')</code>	<code>→ 4</code>
<code>npoints (geometric_type)</code>	<code>npoints (geometric_type) → integer</code> Returns the number of points. Available for path, polygon.	<code>npoints (path ' [(0,0) , (1,1) , (2,0)] ')</code> <code>→ 3</code>
<code>pclose (path)</code>	<code>pclose (path) → path</code> Converts path to closed form.	<code>pclose (path ' [(0,0) , (1,1) , (2,0)] ')</code> <code>→ ((0,0) , (1,1) , (2,0))</code>
<code>popen (path)</code>	<code>popen (path) → path</code> Converts path to open form.	<code>popen (path ' ((0,0) , (1,1) , (2,0)) ')</code> <code>→ [(0,0) , (1,1) , (2,0)]</code>
<code>radius (circle)</code>	<code>radius (circle) → double precision</code> Computes radius of circle.	<code>radius (circle ' < (0,0) , 2 > ')</code> <code>→ 2</code>
<code>slope (point, point)</code>	<code>slope (point, point) → double precision</code> Computes slope of a line drawn through the two points.	<code>slope (point ' (0,0) ' , point ' (2,1) ')</code> <code>→ 0.5</code>
<code>width (box)</code>	<code>width (box) → double precision</code> Computes horizontal size of box.	<code>width (box ' (1,2) , (0,0) ')</code> <code>→ 1</code>

Table 9.38. Geometric Type Conversion Functions

Function	Description	Example(s)
<code>box (circle)</code>	<code>box (circle) → box</code> Computes box inscribed within the circle.	<code>box (circle ' < (0,0) , 2 > ')</code> <code>→ (1.414213562373095, 1.414213562373095) , (</code> <code>-1.414213562373095, -1.414213562373095)</code>
<code>box (point)</code>	<code>box (point) → box</code> Converts point to empty box.	<code>box (point ' (1,0) ')</code> <code>→ (1,0) , (1,0)</code>
<code>box (point, point)</code>	<code>box (point, point) → box</code> Converts any two corner points to box.	<code>box (point ' (0,1) ' , point ' (1,0) ')</code> <code>→ (1,1) , (0,0)</code>
<code>box (polygon)</code>	<code>box (polygon) → box</code> Computes bounding box of polygon.	<code>box (polygon ' ((0,0) , (1,1) , (2,0)) ')</code> <code>→ (2,1) , (0,0)</code>
<code>bound_box (box, box)</code>	<code>bound_box (box, box) → box</code> Computes bounding box of two boxes.	<code>bound_box (box ' (1,1) , (0,0) ' , box ' (4,4) , (3,3) ')</code> <code>→ (4,4) , (0,0)</code>
<code>circle (box)</code>	<code>circle (box) → circle</code> Computes smallest circle enclosing box.	<code>circle (box ' (1,1) , (0,0) ')</code> <code>→ < (0.5,0.5) , 0.7071067811865476 ></code>

Function	Description	Example(s)
<code>circle (point, double precision) → circle</code>	Constructs circle from center and radius.	<code>circle(point ' (0,0) ', 2.0) → <(0,0),2></code>
<code>circle (polygon) → circle</code>	Converts polygon to circle. The circle's center is the mean of the positions of the polygon's points, and the radius is the average distance of the polygon's points from that center.	<code>circle(polygon ' ((0,0), (1,3), (2,0)) ') → <(1,1),1.6094757082487299></code>
<code>line (point, point) → line</code>	Converts two points to the line through them.	<code>line(point ' (-1,0) ', point ' (1,0) ') → {0,-1,0}</code>
<code>lseg (box) → lseg</code>	Extracts box's diagonal as a line segment.	<code>lseg(box ' (1,0), (-1,0) ') → [(1,0), (-1,0)]</code>
<code>lseg (point, point) → lseg</code>	Constructs line segment from two endpoints.	<code>lseg(point ' (-1,0) ', point ' (1,0) ') → [(-1,0), (1,0)]</code>
<code>path (polygon) → path</code>	Converts polygon to a closed path with the same list of points.	<code>path(polygon ' ((0,0), (1,1), (2,0)) ') → ((0,0), (1,1), (2,0))</code>
<code>point (double precision, double precision) → point</code>	Constructs point from its coordinates.	<code>point(23.4, -44.5) → (23.4,-44.5)</code>
<code>point (box) → point</code>	Computes center of box.	<code>point(box ' (1,0), (-1,0) ') → (0,0)</code>
<code>point (circle) → point</code>	Computes center of circle.	<code>point(circle '<(0,0),2>') → (0,0)</code>
<code>point (lseg) → point</code>	Computes center of line segment.	<code>point(lseg ' [(-1,0), (1,0)] ') → (0,0)</code>
<code>point (polygon) → point</code>	Computes center of polygon (the mean of the positions of the polygon's points).	<code>point(polygon ' ((0,0), (1,1), (2,0)) ') → (1,0.3333333333333333)</code>
<code>polygon (box) → polygon</code>	Converts box to a 4-point polygon.	<code>polygon(box ' (1,1), (0,0) ') → ((0,0), (0,1), (1,1), (1,0))</code>
<code>polygon (circle) → polygon</code>	Converts circle to a 12-point polygon.	<code>polygon(circle '<(0,0),2>') → ((-2,0), (-1.7320508075688774, 0.9999999999999999), (-1.0000000000000002, 1.7320508075688772), (-1.2246063538223773e-16, 2), (0.9999999999999996, 1.7320508075688774), (1.732050807568877, 1.0000000000000007), (2, 2.4492127076447545e-16),</code>

Function	Description	Example(s)
		<code>(1.7320508075688776, -0.9999999999999994), (1.0000000000000009, -1.7320508075688767), (3.673819061467132e-16, -2), (-0.9999999999999987, -1.732050807568878), (-1.7320508075688767, -1.0000000000000009))</code>
<code>polygon (integer, circle) → polygon</code>	Converts circle to an <i>n</i> -point polygon.	<code>polygon(4, circle '<(3,0),1>') → ((2,0), (3,1), (4, 1.2246063538223773e-16), (3,-1))</code>
<code>polygon (path) → polygon</code>	Converts closed path to a polygon with the same list of points.	<code>polygon(path '((0,0), (1,1), (2,0))') → ((0,0), (1,1), (2,0))</code>

It is possible to access the two component numbers of a `point` as though the point were an array with indexes 0 and 1. For example, if `t.p` is a point column then `SELECT p[0] FROM t` retrieves the X coordinate and `UPDATE t SET p[1] = ...` changes the Y coordinate. In the same way, a value of type `box` or `lseg` can be treated as an array of two point values.

9.12. Network Address Functions and Operators

The IP network address types, `cidr` and `inet`, support the usual comparison operators shown in [Table 9.1](#) as well as the specialized operators and functions shown in [Table 9.39](#) and [Table 9.40](#).

Any `cidr` value can be cast to `inet` implicitly; therefore, the operators and functions shown below as operating on `inet` also work on `cidr` values. (Where there are separate functions for `inet` and `cidr`, it is because the behavior should be different for the two cases.) Also, it is permitted to cast an `inet` value to `cidr`. When this is done, any bits to the right of the netmask are silently zeroed to create a valid `cidr` value.

Table 9.39. IP Address Operators

Operator	Description	Example(s)
<code>inet << inet → boolean</code>	Is subnet strictly contained by subnet? This operator, and the next four, test for subnet inclusion. They consider only the network parts of the two addresses (ignoring any bits to the right of the netmasks) and determine whether one network is identical to or a subnet of the other.	<code>inet '192.168.1.5' << inet '192.168.1/24' → t</code> <code>inet '192.168.0.5' << inet '192.168.1/24' → f</code> <code>inet '192.168.1/24' << inet '192.168.1/24' → f</code>
<code>inet <=& inet → boolean</code>	Is subnet contained by or equal to subnet?	<code>inet '192.168.1/24' <=& inet '192.168.1/24' → t</code>
<code>inet >> inet → boolean</code>	Does subnet strictly contain subnet?	<code>inet '192.168.1/24' >> inet '192.168.1.5' → t</code>
<code>inet >=& inet → boolean</code>	Does subnet contain or equal subnet?	<code>inet '192.168.1/24' >=& inet '192.168.1/24' → t</code>
<code>inet && inet → boolean</code>		

Operator	Description	Example(s)
	Does either subnet contain or equal the other?	<code>inet '192.168.1/24' && inet '192.168.1.80/28' → t</code> <code>inet '192.168.1/24' && inet '192.168.2.0/28' → f</code>
<code>~ inet → inet</code>	Computes bitwise NOT.	<code>~ inet '192.168.1.6' → 63.87.254.249</code>
<code>inet & inet → inet</code>	Computes bitwise AND.	<code>inet '192.168.1.6' & inet '0.0.0.255' → 0.0.0.6</code>
<code>inet inet → inet</code>	Computes bitwise OR.	<code>inet '192.168.1.6' inet '0.0.0.255' → 192.168.1.255</code>
<code>inet + bigint → inet</code>	Adds an offset to an address.	<code>inet '192.168.1.6' + 25 → 192.168.1.31</code>
<code>bigint + inet → inet</code>	Adds an offset to an address.	<code>200 + inet '::ffff:fff0:1' → ::ffff:255.240.0.201</code>
<code>inet - bigint → inet</code>	Subtracts an offset from an address.	<code>inet '192.168.1.43' - 36 → 192.168.1.7</code>
<code>inet - inet → bigint</code>	Computes the difference of two addresses.	<code>inet '192.168.1.43' - inet '192.168.1.19' → 24</code> <code>inet ':::1' - inet '::ffff:1' → -4294901760</code>

Table 9.40. IP Address Functions

Function	Description	Example(s)
<code>abbrev (inet) → text</code>	Creates an abbreviated display format as text. (The result is the same as the <code>inet</code> output function produces; it is “abbreviated” only in comparison to the result of an explicit cast to <code>text</code> , which for historical reasons will never suppress the netmask part.)	<code>abbrev(inet '10.1.0.0/32') → 10.1.0.0</code>
<code>abbrev (cidr) → text</code>	Creates an abbreviated display format as text. (The abbreviation consists of dropping all-zero octets to the right of the netmask; more examples are in Table 8.22.)	<code>abbrev(cidr '10.1.0.0/16') → 10.1/16</code>
<code>broadcast (inet) → inet</code>	Computes the broadcast address for the address's network.	<code>broadcast(inet '192.168.1.5/24') → 192.168.1.255/24</code>
<code>family (inet) → integer</code>		

Function	Description	Example(s)
	Returns the address's family: 4 for IPv4, 6 for IPv6.	<code>family(inet ':::1') → 6</code>
<code>host (inet) → text</code>	Returns the IP address as text, ignoring the netmask.	<code>host(inet '192.168.1.0/24') → 192.168.1.0</code>
<code>hostmask (inet) → inet</code>	Computes the host mask for the address's network.	<code>hostmask(inet '192.168.23.20/30') → 0.0.0.3</code>
<code>inet_merge (inet, inet) → cidr</code>	Computes the smallest network that includes both of the given networks.	<code>inet_merge(inet '192.168.1.5/24', inet '192.168.2.5/24') → 192.168.0.0/22</code>
<code>inet_same_family (inet, inet) → boolean</code>	Tests whether the addresses belong to the same IP family.	<code>inet_same_family(inet '192.168.1.5/24', inet ':::1') → f</code>
<code>masklen (inet) → integer</code>	Returns the netmask length in bits.	<code>masklen(inet '192.168.1.5/24') → 24</code>
<code>netmask (inet) → inet</code>	Computes the network mask for the address's network.	<code>netmask(inet '192.168.1.5/24') → 255.255.255.0</code>
<code>network (inet) → cidr</code>	Returns the network part of the address, zeroing out whatever is to the right of the netmask. (This is equivalent to casting the value to <code>cidr</code> .)	<code>network(inet '192.168.1.5/24') → 192.168.1.0/24</code>
<code>set_masklen (inet, integer) → inet</code>	Sets the netmask length for an <code>inet</code> value. The address part does not change.	<code>set_masklen(inet '192.168.1.5/24', 16) → 192.168.1.5/16</code>
<code>set_masklen (cidr, integer) → cidr</code>	Sets the netmask length for a <code>cidr</code> value. Address bits to the right of the new netmask are set to zero.	<code>set_masklen(cidr '192.168.1.0/24', 16) → 192.168.0.0/16</code>
<code>text (inet) → text</code>	Returns the unabbreviated IP address and netmask length as text. (This has the same result as an explicit cast to <code>text</code> .)	<code>text(inet '192.168.1.5') → 192.168.1.5/32</code>

Tip

The `abbrev`, `host`, and `text` functions are primarily intended to offer alternative display formats for IP addresses.

The MAC address types, `macaddr` and `macaddr8`, support the usual comparison operators shown in [Table 9.1](#) as well as the specialized functions shown in [Table 9.41](#). In addition, they support the bitwise logical operators `~`, `&` and `|` (NOT, AND and OR), just as shown above for IP addresses.

Table 9.41. MAC Address Functions

Function	Description	Example(s)
<code>trunc (macaddr) → macaddr</code>	Sets the last 3 bytes of the address to zero. The remaining prefix can be associated with a particular manufacturer (using data not included in Postgres Pro).	<code>trunc (macaddr '12:34:56:78:90:ab') → 12:34:56:00:00:00</code>
<code>trunc (macaddr8) → macaddr8</code>	Sets the last 5 bytes of the address to zero. The remaining prefix can be associated with a particular manufacturer (using data not included in Postgres Pro).	<code>trunc (macaddr8 '12:34:56:78:90:ab:cd:ef') → 12:34:56:00:00:00:00:00</code>
<code>macaddr8_set7bit (macaddr8) → macaddr8</code>	Sets the 7th bit of the address to one, creating what is known as modified EUI-64, for inclusion in an IPv6 address.	<code>macaddr8_set7bit (macaddr8 '00:34:56:ab:cd:ef') → 02:34:56:ff:fe:ab:cd:ef</code>

9.13. Text Search Functions and Operators

[Table 9.42](#), [Table 9.43](#) and [Table 9.44](#) summarize the functions and operators that are provided for full text searching. See [Chapter 12](#) for a detailed explanation of Postgres Pro's text search facility.

Table 9.42. Text Search Operators

Operator	Description	Example(s)
<code>tsvector @@ tsquery → boolean</code> <code>tsquery @@ tsvector → boolean</code>	Does <code>tsvector</code> match <code>tsquery</code> ? (The arguments can be given in either order.)	<code>to_tsvector('fat cats ate rats') @@ to_tsquery('cat & rat') → t</code>
<code>text @@ tsquery → boolean</code>	Does text string, after implicit invocation of <code>to_tsvector()</code> , match <code>tsquery</code> ?	<code>'fat cats ate rats' @@ to_tsquery('cat & rat') → t</code>
<code>tsvector tsvector → tsvector</code>	Concatenates two <code>tsvectors</code> . If both inputs contain lexeme positions, the second input's positions are adjusted accordingly.	<code>'a:1 b:2'::tsvector 'c:1 d:2 b:3'::tsvector → 'a':1 'b':2,5 'c':3 'd':4</code>
<code>tsquery && tsquery → tsquery</code>	ANDs two <code>tsqueries</code> together, producing a query that matches documents that match both input queries.	<code>'fat rat'::tsquery && 'cat'::tsquery → ('fat' 'rat') & 'cat'</code>
<code>tsquery tsquery → tsquery</code>	ORs two <code>tsqueries</code> together, producing a query that matches documents that match either input query.	<code>'fat rat'::tsquery 'cat'::tsquery → 'fat' 'rat' 'cat'</code>
<code>!! tsquery → tsquery</code>	Negates a <code>tsquery</code> , producing a query that matches documents that do not match the input query.	<code>!! 'cat'::tsquery → !'cat'</code>

Operator	Description	Example(s)
<code>tsquery <-> tsquery</code>	Constructs a phrase query, which matches if the two input queries match at successive lexemes.	<code>to_tsquery('fat') <-> to_tsquery('rat') → 'fat' <-> 'rat'</code>
<code>tsquery @> tsquery</code>	Does first <code>tsquery</code> contain the second? (This considers only whether all the lexemes appearing in one query appear in the other, ignoring the combining operators.)	<code>'cat'::tsquery @> 'cat & rat'::tsquery → f</code>
<code>tsquery <@ tsquery</code>	Is first <code>tsquery</code> contained in the second? (This considers only whether all the lexemes appearing in one query appear in the other, ignoring the combining operators.)	<code>'cat'::tsquery <@ 'cat & rat'::tsquery → t</code> <code>'cat'::tsquery <@ '!cat & rat'::tsquery → t</code>

In addition to these specialized operators, the usual comparison operators shown in [Table 9.1](#) are available for types `tsvector` and `tsquery`. These are not very useful for text searching but allow, for example, unique indexes to be built on columns of these types.

Table 9.43. Text Search Functions

Function	Description	Example(s)
<code>array_to_tsvector (text[])</code>	Converts an array of text strings to a <code>tsvector</code> . The given strings are used as lexemes as-is, without further processing. Array elements must not be empty strings or <code>NULL</code> .	<code>array_to_tsvector('{fat,cat,rat'}::text[]) → 'cat' 'fat' 'rat'</code>
<code>get_current_ts_config ()</code>	Returns the OID of the current default text search configuration (as set by default_text_search_config).	<code>get_current_ts_config() → english</code>
<code>length (tsvector)</code>	Returns the number of lexemes in the <code>tsvector</code> .	<code>length('fat:2,4 cat:3 rat:5A'::tsvector) → 3</code>
<code>numnode (tsquery)</code>	Returns the number of lexemes plus operators in the <code>tsquery</code> .	<code>numnode('(fat & rat) cat'::tsquery) → 5</code>
<code>plainto_tsquery ([config regconfig,] query text)</code>	Converts text to a <code>tsquery</code> , normalizing words according to the specified or default configuration. Any punctuation in the string is ignored (it does not determine query operators). The resulting query matches documents containing all non-stopwords in the text.	<code>plainto_tsquery('english', 'The Fat Rats') → 'fat' & 'rat'</code>
<code>phraseto_tsquery ([config regconfig,] query text)</code>	Converts text to a <code>tsquery</code> , normalizing words according to the specified or default configuration. Any punctuation in the string is ignored (it does not determine query operators). The resulting query matches phrases containing all non-stopwords in the text.	<code>phraseto_tsquery('english', 'The Fat Rats') → 'fat' <-> 'rat'</code>

Function	Description	Example(s)
<code>phraseto_tsquery</code>		<code>phraseto_tsquery('english', 'The Cat and Rats') → 'cat' <2> 'rat'</code>
<code>websearch_to_tsquery</code>	Converts text to a <code>tsquery</code> , normalizing words according to the specified or default configuration. Quoted word sequences are converted to phrase tests. The word “or” is understood as producing an OR operator, and a dash produces a NOT operator; other punctuation is ignored. This approximates the behavior of some common web search tools.	<code>websearch_to_tsquery('english', '"fat rat" or cat dog') → 'fat' <-> 'rat' 'cat' & 'dog'</code>
<code>querytree</code>	Produces a representation of the indexable portion of a <code>tsquery</code> . A result that is empty or just <code>T</code> indicates a non-indexable query.	<code>querytree('foo & ! bar'::tsquery) → 'foo'</code>
<code>setweight</code>	Assigns the specified <i>weight</i> to each element of the <i>vector</i> .	<code>setweight('fat:2,4 cat:3 rat:5B'::tsvector, 'A') → 'cat':3A 'fat':2A,4A 'rat':5A</code>
<code>setweight</code>	Assigns the specified <i>weight</i> to elements of the <i>vector</i> that are listed in <i>lexemes</i> . The strings in <i>lexemes</i> are taken as lexemes as-is, without further processing. Strings that do not match any lexeme in <i>vector</i> are ignored.	<code>setweight('fat:2,4 cat:3 rat:5,6B'::tsvector, 'A', '{cat,rat}') → 'cat':3A 'fat':2,4 'rat':5A,6A</code>
<code>strip</code>	Removes positions and weights from the <code>tsvector</code> .	<code>strip('fat:2,4 cat:3 rat:5A'::tsvector) → 'cat' 'fat' 'rat'</code>
<code>to_tsquery</code>	Converts text to a <code>tsquery</code> , normalizing words according to the specified or default configuration. The words must be combined by valid <code>tsquery</code> operators.	<code>to_tsquery('english', 'The & Fat & Rats') → 'fat' & 'rat'</code>
<code>to_tsvector</code>	Converts text to a <code>tsvector</code> , normalizing words according to the specified or default configuration. Position information is included in the result.	<code>to_tsvector('english', 'The Fat Rats') → 'fat':2 'rat':3</code>
<code>to_tsvector</code>	Converts each string value in the JSON document to a <code>tsvector</code> , normalizing words according to the specified or default configuration. The results are then concatenated in document order to produce the output. Position information is generated as though one stopword exists between each pair of string values. (Beware that “document order” of the fields of a JSON object is implementation-dependent when the input is <code>jsonb</code> ; observe the difference in the examples.)	<code>to_tsvector('english', '{"aa": "The Fat Rats", "b": "dog"}'::json) → 'dog':5 'fat':2 'rat':3</code> <code>to_tsvector('english', '{"aa": "The Fat Rats", "b": "dog"}'::jsonb) → 'dog':1 'fat':4 'rat':5</code>

Function	Description	Example(s)
<code>json_to_tsvector ([config regconfig,] document json, filter jsonb) → tsvector</code> <code>jsonb_to_tsvector ([config regconfig,] document jsonb, filter jsonb) → tsvector</code>	<p>Selects each item in the JSON document that is requested by the <i>filter</i> and converts each one to a <i>tsvector</i>, normalizing words according to the specified or default configuration. The results are then concatenated in document order to produce the output. Position information is generated as though one stopword exists between each pair of selected items. (Beware that “document order” of the fields of a JSON object is implementation-dependent when the input is <i>jsonb</i>.) The <i>filter</i> must be a <i>jsonb</i> array containing zero or more of these keywords: “string” (to include all string values), “numeric” (to include all numeric values), “boolean” (to include all boolean values), “key” (to include all keys), or “all” (to include all the above). As a special case, the <i>filter</i> can also be a simple JSON value that is one of these keywords.</p>	<pre>json_to_tsvector('english', '{"a": "The Fat Rats", "b": 123}':::json, '["string", "numeric"]') → '123':5 'fat':2 'rat':3 json_to_tsvector('english', '{"cat": "The Fat Rats", "dog": 123}':::json, 'all') → '123':9 'cat':1 'dog':7 'fat':4 'rat':5</pre>
<code>ts_delete (vector tsvector, lexeme text) → tsvector</code>	<p>Removes any occurrence of the given <i>lexeme</i> from the <i>vector</i>. The <i>lexeme</i> string is treated as a lexeme as-is, without further processing.</p>	<pre>ts_delete('fat:2,4 cat:3 rat:5A':::tsvector, 'fat') → 'cat':3 'rat':5A</pre>
<code>ts_delete (vector tsvector, lexemes text[]) → tsvector</code>	<p>Removes any occurrences of the lexemes in <i>lexemes</i> from the <i>vector</i>. The strings in <i>lexemes</i> are taken as lexemes as-is, without further processing. Strings that do not match any lexeme in <i>vector</i> are ignored.</p>	<pre>ts_delete('fat:2,4 cat:3 rat:5A':::tsvector, ARRAY['fat','rat']) → 'cat':3</pre>
<code>ts_filter (vector tsvector, weights "char"[]) → tsvector</code>	<p>Selects only elements with the given <i>weights</i> from the <i>vector</i>.</p>	<pre>ts_filter('fat:2,4 cat:3b,7c rat:5A':::tsvector, '{a,b}') → 'cat':3B 'rat':5A</pre>
<code>ts_headline ([config regconfig,] document text, query tsquery [, options text]) → text</code>	<p>Displays, in an abbreviated form, the match(es) for the <i>query</i> in the <i>document</i>, which must be raw text not a <i>tsvector</i>. Words in the document are normalized according to the specified or default configuration before matching to the query. Use of this function is discussed in Section 12.3.4, which also describes the available <i>options</i>.</p>	<pre>ts_headline('The fat cat ate the rat.', 'cat') → The fat cat ate the rat.</pre>
<code>ts_headline ([config regconfig,] document json, query tsquery [, options text]) → text</code> <code>ts_headline ([config regconfig,] document jsonb, query tsquery [, options text]) → text</code>	<p>Displays, in an abbreviated form, match(es) for the <i>query</i> that occur in string values within the JSON <i>document</i>. See Section 12.3.4 for more details.</p>	<pre>ts_headline('{"cat":"raining cats and dogs}':::jsonb, 'cat') → {"cat": "rain- ing cats and dogs"}</pre>
<code>ts_rank ([weights real[],] vector tsvector, query tsquery [, normalization integer]) → real</code>	<p>Computes a score showing how well the <i>vector</i> matches the <i>query</i>. See Section 12.3.3 for details.</p>	<pre>ts_rank(to_tsvector('raining cats and dogs'), 'cat') → 0.06079271</pre>

Function	Description	Example(s)												
<code>ts_rank_cd</code>	([<i>weights</i> real[],] <i>vector</i> tsvector, <i>query</i> tsquery [, <i>normalization</i> integer]) → real Computes a score showing how well the <i>vector</i> matches the <i>query</i> , using a cover density algorithm. See Section 12.3.3 for details.	<code>ts_rank_cd(to_tsvector('raining cats and dogs'), 'cat') → 0.1</code>												
<code>ts_rewrite</code>	(<i>query</i> tsquery, <i>target</i> tsquery, <i>substitute</i> tsquery) → tsquery Replaces occurrences of <i>target</i> with <i>substitute</i> within the <i>query</i> . See Section 12.4.2.1 for details.	<code>ts_rewrite('a & b'::tsquery, 'a'::tsquery, 'foo bar'::tsquery) → 'b' & ('foo' 'bar')</code>												
<code>ts_rewrite</code>	(<i>query</i> tsquery, <i>select</i> text) → tsquery Replaces portions of the <i>query</i> according to <i>target</i> (s) and <i>substitute</i> (s) obtained by executing a SELECT command. See Section 12.4.2.1 for details.	<code>SELECT ts_rewrite('a & b'::tsquery, 'SELECT t,s FROM aliases') → 'b' & ('foo' 'bar')</code>												
<code>tsquery_phrase</code>	(<i>query1</i> tsquery, <i>query2</i> tsquery) → tsquery Constructs a phrase query that searches for matches of <i>query1</i> and <i>query2</i> at successive lexemes (same as <-> operator).	<code>tsquery_phrase(to_tsquery('fat'), to_tsquery('cat')) → 'fat' <-> 'cat'</code>												
<code>tsquery_phrase</code>	(<i>query1</i> tsquery, <i>query2</i> tsquery, <i>distance</i> integer) → tsquery Constructs a phrase query that searches for matches of <i>query1</i> and <i>query2</i> that occur exactly <i>distance</i> lexemes apart.	<code>tsquery_phrase(to_tsquery('fat'), to_tsquery('cat'), 10) → 'fat' <10> 'cat'</code>												
<code>tsvector_to_array</code>	(tsvector) → text[] Converts a tsvector to an array of lexemes.	<code>tsvector_to_array('fat:2,4 cat:3 rat:5A'::tsvector) → {cat,fat,rat}</code>												
<code>unnest</code>	(tsvector) → setof record (<i>lexeme</i> text, <i>positions</i> smallint[], <i>weights</i> text) Expands a tsvector into a set of rows, one per lexeme.	<code>select * from unnest('cat:3 fat:2,4 rat:5A'::tsvector) →</code> <table><thead><tr><th>lexeme</th><th>positions</th><th>weights</th></tr></thead><tbody><tr><td>cat</td><td>{3}</td><td>{D}</td></tr><tr><td>fat</td><td>{2,4}</td><td>{D,D}</td></tr><tr><td>rat</td><td>{5}</td><td>{A}</td></tr></tbody></table>	lexeme	positions	weights	cat	{3}	{D}	fat	{2,4}	{D,D}	rat	{5}	{A}
lexeme	positions	weights												
cat	{3}	{D}												
fat	{2,4}	{D,D}												
rat	{5}	{A}												

Note

All the text search functions that accept an optional `regconfig` argument will use the configuration specified by [default_text_search_config](#) when that argument is omitted.

The functions in [Table 9.44](#) are listed separately because they are not usually used in everyday text searching operations. They are primarily helpful for development and debugging of new text search configurations.

Table 9.44. Text Search Debugging Functions

Function	Description	Example(s)
<code>ts_debug</code>	<code>([config regconfig,] document text) → setof record (alias text, description text, token text, dictionaries regdictionary[], dictionary regdictionary, lexemes text[])</code> Extracts and normalizes tokens from the <i>document</i> according to the specified or default text search configuration, and returns information about how each token was processed. See Section 12.8.1 for details.	<code>ts_debug('english', 'The Brightest supernovaes') → (asciiword, "Word, all ASCII", The, {english_stem}, english_stem, { }) ...</code>
<code>ts_lexize</code>	<code>(dict regdictionary, token text) → text[]</code> Returns an array of replacement lexemes if the input token is known to the dictionary, or an empty array if the token is known to the dictionary but it is a stop word, or NULL if it is not a known word. See Section 12.8.3 for details.	<code>ts_lexize('english_stem', 'stars') → {star}</code>
<code>ts_parse</code>	<code>(parser_name text, document text) → setof record (tokid integer, token text)</code> Extracts tokens from the <i>document</i> using the named parser. See Section 12.8.2 for details.	<code>ts_parse('default', 'foo - bar') → (1,foo) ...</code>
<code>ts_parse</code>	<code>(parser_oid oid, document text) → setof record (tokid integer, token text)</code> Extracts tokens from the <i>document</i> using a parser specified by OID. See Section 12.8.2 for details.	<code>ts_parse(3722, 'foo - bar') → (1,foo) ...</code>
<code>ts_token_type</code>	<code>(parser_name text) → setof record (tokid integer, alias text, description text)</code> Returns a table that describes each type of token the named parser can recognize. See Section 12.8.2 for details.	<code>ts_token_type('default') → (1,asciiword, "Word, all ASCII") ...</code>
<code>ts_token_type</code>	<code>(parser_oid oid) → setof record (tokid integer, alias text, description text)</code> Returns a table that describes each type of token a parser specified by OID can recognize. See Section 12.8.2 for details.	<code>ts_token_type(3722) → (1,asciiword, "Word, all ASCII") ...</code>
<code>ts_stat</code>	<code>(sqlquery text [, weights text]) → setof record (word text, ndoc integer, nentry integer)</code> Executes the <i>sqlquery</i> , which must return a single <i>tsvector</i> column, and returns statistics about each distinct lexeme contained in the data. See Section 12.4.4 for details.	<code>ts_stat('SELECT vector FROM apod') → (foo,10,15) ...</code>

9.14. UUID Functions

Postgres Pro includes one function to generate a UUID:

`gen_random_uuid () → uuid`

This function returns a version 4 (random) UUID. This is the most commonly used type of UUID and is appropriate for most applications.

The [uuid-oss](#) module provides additional functions that implement other standard algorithms for generating UUIDs.

There are also functions to extract data from UUIDs:

`uuid_extract_timestamp (uuid) → timestamp with time zone`

This function extracts a timestamp with time zone from UUID version 1. For other versions, this function returns null. Note that the extracted timestamp is not necessarily exactly equal to the time the UUID was generated; this depends on the implementation that generated the UUID.

`uuid_extract_version (uuid) → smallint`

This function extracts the version from a UUID of the variant described by [RFC 4122](#). For other variants, this function returns null. For example, for a UUID generated by `gen_random_uuid`, this function will return 4.

Postgres Pro also provides the usual comparison operators shown in [Table 9.1](#) for UUIDs.

9.15. XML Functions

The functions and function-like expressions described in this section operate on values of type `xml`. See [Section 8.13](#) for information about the `xml` type. The function-like expressions `xmlparse` and `xmlserialize` for converting to and from type `xml` are documented there, not in this section.

Use of most of these functions requires Postgres Pro to have been built with `configure --with-libxml`.

9.15.1. Producing XML Content

A set of functions and function-like expressions is available for producing XML content from SQL data. As such, they are particularly suitable for formatting query results into XML documents for processing in client applications.

9.15.1.1. `xmltext`

`xmltext (text) → xml`

The function `xmltext` returns an XML value with a single text node containing the input argument as its content. Predefined entities like ampersand (&), left and right angle brackets (< >), and quotation marks (") are escaped.

Example:

```
SELECT xmltext('< foo & bar >');
      xmltext
-----
<lt; foo & bar >gt;
```

9.15.1.2. `xmlcomment`

`xmlcomment (text) → xml`

The function `xmlcomment` creates an XML value containing an XML comment with the specified text as content. The text cannot contain “--” or end with a “-”, otherwise the resulting construct would not be a valid XML comment. If the argument is null, the result is null.

Example:

```
SELECT xmlcomment('hello');
      xmlcomment
-----
<!--hello-->
```

9.15.1.3. `xmlconcat`

`xmlconcat (xml [, ...]) → xml`

The function `xmlconcat` concatenates a list of individual XML values to create a single value containing an XML content fragment. Null values are omitted; the result is only null if there are no nonnull arguments.

Example:

```
SELECT xmlconcat ('<abc/>', '<bar>foo</bar>');
```

```
      xmlconcat
-----
<abc/><bar>foo</bar>
```

XML declarations, if present, are combined as follows. If all argument values have the same XML version declaration, that version is used in the result, else no version is used. If all argument values have the standalone declaration value “yes”, then that value is used in the result. If all argument values have a standalone declaration value and at least one is “no”, then that is used in the result. Else the result will have no standalone declaration. If the result is determined to require a standalone declaration but no version declaration, a version declaration with version 1.0 will be used because XML requires an XML declaration to contain a version declaration. Encoding declarations are ignored and removed in all cases.

Example:

```
SELECT xmlconcat ('<?xml version="1.1"?><foo/>', '<?xml version="1.1" standalone="no"?><bar/>');
```

```
      xmlconcat
-----
<?xml version="1.1"?><foo/><bar/>
```

9.15.1.4. xmlelement

```
xmlelement ( NAME name [, XMLATTRIBUTES ( attvalue [ AS attname ] [, ...] ) ]
            [, content [, ...]] ) → xml
```

The `xmlelement` expression produces an XML element with the given name, attributes, and content. The *name* and *attname* items shown in the syntax are simple identifiers, not values. The *attvalue* and *content* items are expressions, which can yield any Postgres Pro data type. The argument(s) within XMLATTRIBUTES generate attributes of the XML element; the *content* value(s) are concatenated to form its content.

Examples:

```
SELECT xmlelement(name foo);
```

```
      xmlelement
-----
<foo/>
```

```
SELECT xmlelement(name foo, xmlattributes('xyz' as bar));
```

```
      xmlelement
-----
<foo bar="xyz"/>
```

```
SELECT xmlelement(name foo, xmlattributes(current_date as bar), 'cont', 'ent');
```

```
      xmlelement
-----
<foo bar="2007-01-26">content</foo>
```

Element and attribute names that are not valid XML names are escaped by replacing the offending characters by the sequence `_xHHHH_`, where *HHHH* is the character's Unicode codepoint in hexadecimal notation. For example:

```
SELECT xmlelement(name "foo$bar", xmlattributes('xyz' as "a&b"));
```

xmlelement

```
-----
<foo_x0024_bar a_x0026_b="xyz"/>
```

An explicit attribute name need not be specified if the attribute value is a column reference, in which case the column's name will be used as the attribute name by default. In other cases, the attribute must be given an explicit name. So this example is valid:

```
CREATE TABLE test (a xml, b xml);
SELECT xmlelement(name test, xmlattributes(a, b)) FROM test;
```

But these are not:

```
SELECT xmlelement(name test, xmlattributes('constant'), a, b) FROM test;
SELECT xmlelement(name test, xmlattributes(func(a, b))) FROM test;
```

Element content, if specified, will be formatted according to its data type. If the content is itself of type xml, complex XML documents can be constructed. For example:

```
SELECT xmlelement(name foo, xmlattributes('xyz' as bar),
                  xmlelement(name abc),
                  xmlcomment('test'),
                  xmlelement(name xyz));
```

xmlelement

```
-----
<foo bar="xyz"><abc/><!--test--><xyz/></foo>
```

Content of other types will be formatted into valid XML character data. This means in particular that the characters <, >, and & will be converted to entities. Binary data (data type `bytea`) will be represented in base64 or hex encoding, depending on the setting of the configuration parameter [xmlbinary](#). The particular behavior for individual data types is expected to evolve in order to align the Postgres Pro mappings with those specified in SQL:2006 and later, as discussed in [Section D.3.1.3](#).

9.15.1.5. xmlforest

```
xmlforest ( content [ AS name ] [, ...] ) → xml
```

The `xmlforest` expression produces an XML forest (sequence) of elements using the given names and content. As for `xmlelement`, each *name* must be a simple identifier, while the *content* expressions can have any data type.

Examples:

```
SELECT xmlforest('abc' AS foo, 123 AS bar);
```

xmlforest

```
-----
<foo>abc</foo><bar>123</bar>
```

```
SELECT xmlforest(table_name, column_name)
FROM information_schema.columns
WHERE table_schema = 'pg_catalog';
```

xmlforest

```
-----
<table_name>pg_authid</table_name><column_name>rolname</column_name>
<table_name>pg_authid</table_name><column_name>rolsuper</column_name>
...
```

As seen in the second example, the element name can be omitted if the content value is a column reference, in which case the column name is used by default. Otherwise, a name must be specified.

Element names that are not valid XML names are escaped as shown for `xmlelement` above. Similarly, content data is escaped to make valid XML content, unless it is already of type `xml`.

Note that XML forests are not valid XML documents if they consist of more than one element, so it might be useful to wrap `xmlforest` expressions in `xmlelement`.

9.15.1.6. `xmlpi`

```
xmlpi ( NAME name [, content ] ) → xml
```

The `xmlpi` expression creates an XML processing instruction. As for `xmlelement`, the `name` must be a simple identifier, while the `content` expression can have any data type. The `content`, if present, must not contain the character sequence `>`.

Example:

```
SELECT xmlpi(name php, 'echo "hello world";');

      xmlpi
-----
<?php echo "hello world";?>
```

9.15.1.7. `xmlroot`

```
xmlroot ( xml, VERSION {text|NO VALUE} [, STANDALONE {YES|NO|NO VALUE} ] ) → xml
```

The `xmlroot` expression alters the properties of the root node of an XML value. If a version is specified, it replaces the value in the root node's version declaration; if a standalone setting is specified, it replaces the value in the root node's standalone declaration.

```
SELECT xmlroot(xmlparse(document '<?xml version="1.1"?><content>abc</content>'),
              version '1.0', standalone yes);

      xmlroot
-----
<?xml version="1.0" standalone="yes"?>
<content>abc</content>
```

9.15.1.8. `xmlagg`

```
xmlagg ( xml ) → xml
```

The function `xmlagg` is, unlike the other functions described here, an aggregate function. It concatenates the input values to the aggregate function call, much like `xmlconcat` does, except that concatenation occurs across rows rather than across expressions in a single row. See [Section 9.21](#) for additional information about aggregate functions.

Example:

```
CREATE TABLE test (y int, x xml);
INSERT INTO test VALUES (1, '<foo>abc</foo>');
INSERT INTO test VALUES (2, '<bar/>');
SELECT xmlagg(x) FROM test;

      xmlagg
-----
<foo>abc</foo><bar/>
```

To determine the order of the concatenation, an `ORDER BY` clause may be added to the aggregate call as described in [Section 4.2.7](#). For example:

```
SELECT xmlagg(x ORDER BY y DESC) FROM test;

      xmlagg
-----
<bar/><foo>abc</foo>
```

The following non-standard approach used to be recommended in previous versions, and may still be useful in specific cases:

```
SELECT xmlagg(x) FROM (SELECT * FROM test ORDER BY y DESC) AS tab;
      xmlagg
-----
<bar/><foo>abc</foo>
```

9.15.2. XML Predicates

The expressions described in this section check properties of `xml` values.

9.15.2.1. IS DOCUMENT

```
xml IS DOCUMENT → boolean
```

The expression `IS DOCUMENT` returns true if the argument XML value is a proper XML document, false if it is not (that is, it is a content fragment), or null if the argument is null. See [Section 8.13](#) about the difference between documents and content fragments.

9.15.2.2. IS NOT DOCUMENT

```
xml IS NOT DOCUMENT → boolean
```

The expression `IS NOT DOCUMENT` returns false if the argument XML value is a proper XML document, true if it is not (that is, it is a content fragment), or null if the argument is null.

9.15.2.3. XMLEXISTS

```
XMLEXISTS ( text PASSING [BY {REF|VALUE}] xml [BY {REF|VALUE}] ) → boolean
```

The function `xmlexists` evaluates an XPath 1.0 expression (the first argument), with the passed XML value as its context item. The function returns false if the result of that evaluation yields an empty node-set, true if it yields any other value. The function returns null if any argument is null. A nonnull value passed as the context item must be an XML document, not a content fragment or any non-XML value.

Example:

```
SELECT xmlexists('//town[text() = ''Toronto'']' PASSING BY VALUE
  '<towns><town>Toronto</town><town>Ottawa</town></towns>');

 xmlexists
-----
 t
(1 row)
```

The `BY REF` and `BY VALUE` clauses are accepted in Postgres Pro, but are ignored, as discussed in [Section D.3.2](#).

In the SQL standard, the `xmlexists` function evaluates an expression in the XML Query language, but Postgres Pro allows only an XPath 1.0 expression, as discussed in [Section D.3.1](#).

9.15.2.4. xml_is_well_formed

```
xml_is_well_formed ( text ) → boolean
xml_is_well_formed_document ( text ) → boolean
xml_is_well_formed_content ( text ) → boolean
```

These functions check whether a text string represents well-formed XML, returning a Boolean result. `xml_is_well_formed_document` checks for a well-formed document, while `xml_is_well_formed_content` checks for well-formed content. `xml_is_well_formed` does the former if the `xmloption` configuration parameter is set to `DOCUMENT`, or the latter if it is set to `CONTENT`. This means that `xml_is_well_formed` is useful for seeing whether a simple cast to type `xml` will succeed, whereas the other two functions are useful for seeing whether the corresponding variants of `XMLPARSE` will succeed.

Examples:

```
SET xmloption TO DOCUMENT;
SELECT xml_is_well_formed('<>');
xml_is_well_formed
```

```
-----
f
(1 row)
```

```
SELECT xml_is_well_formed('<abc/>');
xml_is_well_formed
```

```
-----
t
(1 row)
```

```
SET xmloption TO CONTENT;
SELECT xml_is_well_formed('abc');
xml_is_well_formed
```

```
-----
t
(1 row)
```

```
SELECT xml_is_well_formed_document('<pg:foo xmlns:pg="http://postgresql.org/
stuff">bar</pg:foo>');
xml_is_well_formed_document
```

```
-----
t
(1 row)
```

```
SELECT xml_is_well_formed_document('<pg:foo xmlns:pg="http://postgresql.org/
stuff">bar</my:foo>');
xml_is_well_formed_document
```

```
-----
f
(1 row)
```

The last example shows that the checks include whether namespaces are correctly matched.

9.15.3. Processing XML

To process values of data type `xml`, Postgres Pro offers the functions `xpath` and `xpath_exists`, which evaluate XPath 1.0 expressions, and the `XMLTABLE` table function.

9.15.3.1. `xpath`

```
xpath ( xpath text, xml xml [, nsarray text[] ] ) → xml[]
```

The function `xpath` evaluates the XPath 1.0 expression `xpath` (given as `text`) against the XML value `xml`. It returns an array of XML values corresponding to the node-set produced by the XPath expression. If the XPath expression returns a scalar value rather than a node-set, a single-element array is returned.

The second argument must be a well formed XML document. In particular, it must have a single root node element.

The optional third argument of the function is an array of namespace mappings. This array should be a two-dimensional `text` array with the length of the second axis being equal to 2 (i.e., it should be an array of arrays, each of which consists of exactly 2 elements). The first element of each array entry is the namespace name (alias), the second the namespace URI. It is not required that aliases provided in this array be the same as those being used in the XML document itself (in other words, both in the XML document and in the `xpath` function context, aliases are *local*).

Example:

```
SELECT xpath('/my:a/text()', '<my:a xmlns:my="http://example.com">test</my:a>',
           ARRAY[ARRAY['my', 'http://example.com']]);
```

```
xpath
-----
{test}
(1 row)
```

To deal with default (anonymous) namespaces, do something like this:

```
SELECT xpath('//mydefns:b/text()', '<a xmlns="http://example.com"><b>test</b></a>',
           ARRAY[ARRAY['mydefns', 'http://example.com']]);
```

```
xpath
-----
{test}
(1 row)
```

9.15.3.2. xpath_exists

```
xpath_exists ( xpath text, xml xml [, nsarray text[] ] ) → boolean
```

The function `xpath_exists` is a specialized form of the `xpath` function. Instead of returning the individual XML values that satisfy the XPath 1.0 expression, this function returns a Boolean indicating whether the query was satisfied or not (specifically, whether it produced any value other than an empty node-set). This function is equivalent to the `XMLEXISTS` predicate, except that it also offers support for a namespace mapping argument.

Example:

```
SELECT xpath_exists('/my:a/text()', '<my:a xmlns:my="http://example.com">test</my:a>',
                   ARRAY[ARRAY['my', 'http://example.com']]);
```

```
xpath_exists
-----
t
(1 row)
```

9.15.3.3. xmltable

```
XMLTABLE (
  [ XMLNAMESPACES ( namespace_uri AS namespace_name [, ...] ), ]
  row_expression PASSING [BY {REF|VALUE}] document_expression [BY {REF|VALUE}]
  COLUMNS name { type [PATH column_expression] [DEFAULT default_expression] [NOT NULL
| NULL]
               | FOR ORDINALITY }
  [, ...]
) → setof record
```

The `xmltable` expression produces a table based on an XML value, an XPath filter to extract rows, and a set of column definitions. Although it syntactically resembles a function, it can only appear as a table in a query's `FROM` clause.

The optional `XMLNAMESPACES` clause gives a comma-separated list of namespace definitions, where each `namespace_uri` is a text expression and each `namespace_name` is a simple identifier. It specifies the XML namespaces used in the document and their aliases. A default namespace specification is not currently supported.

The required `row_expression` argument is an XPath 1.0 expression (given as `text`) that is evaluated, passing the XML value `document_expression` as its context item, to obtain a set of XML nodes. These

nodes are what `xmltable` transforms into output rows. No rows will be produced if the *document_expression* is null, nor if the *row_expression* produces an empty node-set or any value other than a node-set.

document_expression provides the context item for the *row_expression*. It must be a well-formed XML document; fragments/forests are not accepted. The `BY REF` and `BY VALUE` clauses are accepted but ignored, as discussed in [Section D.3.2](#).

In the SQL standard, the `xmltable` function evaluates expressions in the XML Query language, but Postgres Pro allows only XPath 1.0 expressions, as discussed in [Section D.3.1](#).

The required `COLUMNS` clause specifies the column(s) that will be produced in the output table. See the syntax summary above for the format. A name is required for each column, as is a data type (unless `FOR ORDINALITY` is specified, in which case type `integer` is implicit). The path, default and nullability clauses are optional.

A column marked `FOR ORDINALITY` will be populated with row numbers, starting with 1, in the order of nodes retrieved from the *row_expression*'s result node-set. At most one column may be marked `FOR ORDINALITY`.

Note

XPath 1.0 does not specify an order for nodes in a node-set, so code that relies on a particular order of the results will be implementation-dependent. Details can be found in [Section D.3.1.2](#).

The *column_expression* for a column is an XPath 1.0 expression that is evaluated for each row, with the current node from the *row_expression* result as its context item, to find the value of the column. If no *column_expression* is given, then the column name is used as an implicit path.

If a column's XPath expression returns a non-XML value (which is limited to string, boolean, or double in XPath 1.0) and the column has a Postgres Pro type other than `xml`, the column will be set as if by assigning the value's string representation to the Postgres Pro type. (If the value is a boolean, its string representation is taken to be 1 or 0 if the output column's type category is numeric, otherwise `true` or `false`.)

If a column's XPath expression returns a non-empty set of XML nodes and the column's Postgres Pro type is `xml`, the column will be assigned the expression result exactly, if it is of document or content form.¹

A non-XML result assigned to an `xml` output column produces content, a single text node with the string value of the result. An XML result assigned to a column of any other type may not have more than one node, or an error is raised. If there is exactly one node, the column will be set as if by assigning the node's string value (as defined for the XPath 1.0 `string` function) to the Postgres Pro type.

The string value of an XML element is the concatenation, in document order, of all text nodes contained in that element and its descendants. The string value of an element with no descendant text nodes is an empty string (not `NULL`). Any `xsi:nil` attributes are ignored. Note that the whitespace-only `text()` node between two non-text elements is preserved, and that leading whitespace on a `text()` node is not flattened. The XPath 1.0 `string` function may be consulted for the rules defining the string value of other XML node types and non-XML values.

The conversion rules presented here are not exactly those of the SQL standard, as discussed in [Section D.3.1.3](#).

If the path expression returns an empty node-set (typically, when it does not match) for a given row, the column will be set to `NULL`, unless a *default_expression* is specified; then the value resulting from evaluating that expression is used.

¹ A result containing more than one element node at the top level, or non-whitespace text outside of an element, is an example of content form. An XPath result can be of neither form, for example if it returns an attribute node selected from the element that contains it. Such a result will be put into content form with each such disallowed node replaced by its string value, as defined for the XPath 1.0 `string` function.

A *default_expression*, rather than being evaluated immediately when `xmltable` is called, is evaluated each time a default is needed for the column. If the expression qualifies as stable or immutable, the repeat evaluation may be skipped. This means that you can usefully use volatile functions like `nextval` in *default_expression*.

Columns may be marked `NOT NULL`. If the *column_expression* for a `NOT NULL` column does not match anything and there is no `DEFAULT` or the *default_expression* also evaluates to null, an error is reported.

Examples:

```
CREATE TABLE xmldata AS SELECT
xml $$
<ROWS>
  <ROW id="1">
    <COUNTRY_ID>AU</COUNTRY_ID>
    <COUNTRY_NAME>Australia</COUNTRY_NAME>
  </ROW>
  <ROW id="5">
    <COUNTRY_ID>JP</COUNTRY_ID>
    <COUNTRY_NAME>Japan</COUNTRY_NAME>
    <PREMIER_NAME>Shinzo Abe</PREMIER_NAME>
    <SIZE unit="sq_mi">145935</SIZE>
  </ROW>
  <ROW id="6">
    <COUNTRY_ID>SG</COUNTRY_ID>
    <COUNTRY_NAME>Singapore</COUNTRY_NAME>
    <SIZE unit="sq_km">697</SIZE>
  </ROW>
</ROWS>
$$ AS data;

SELECT xmltable.*
FROM xmldata,
     XMLTABLE ('//ROWS/ROW'
               PASSING data
               COLUMNS id int PATH '@id',
                        ordinality FOR ORDINALITY,
                        "COUNTRY_NAME" text,
                        country_id text PATH 'COUNTRY_ID',
                        size_sq_km float PATH 'SIZE[@unit = "sq_km"]',
                        size_other text PATH
                          'concat(SIZE[@unit!="sq_km"], " ", SIZE[@unit!="sq_km"]/'
@unit) ',
                        premier_name text PATH 'PREMIER_NAME' DEFAULT 'not specified');

 id | ordinality | COUNTRY_NAME | country_id | size_sq_km | size_other | premier_name
----+-----+-----+-----+-----+-----+-----
 1 |          1 | Australia   | AU         |             |             | not
specified
 5 |          2 | Japan       | JP         |             | 145935 sq_mi | Shinzo Abe
 6 |          3 | Singapore   | SG         | 697         |             | not
specified
```

The following example shows concatenation of multiple `text()` nodes, usage of the column name as XPath filter, and the treatment of whitespace, XML comments and processing instructions:

```
CREATE TABLE xmlelements AS SELECT
xml $$
  <root>
```



```

    <element> Hello<!-- xyxxz -->2a2<?aaaaa?> <!--x--> bbb<x>xxx</x>CC </element>
</root>
$$ AS data;

SELECT xmltable.*
FROM xmlelements, XMLTABLE('/root' PASSING data COLUMNS element text);
      element
-----
Hello2a2   bbbxxxCC

```

The following example illustrates how the `XMLNAMESPACES` clause can be used to specify a list of namespaces used in the XML document as well as in the XPath expressions:

```

WITH xmldata(data) AS (VALUES ('
<example xmlns="http://example.com/myns" xmlns:B="http://example.com/b">
  <item foo="1" B:bar="2"/>
  <item foo="3" B:bar="4"/>
  <item foo="4" B:bar="5"/>
</example>'::xml)
)
SELECT xmltable.*
FROM XMLTABLE(XMLNAMESPACES('http://example.com/myns' AS x,
                           'http://example.com/b' AS "B"),
              '/x:example/x:item'
              PASSING (SELECT data FROM xmldata)
              COLUMNS foo int PATH '@foo',
                       bar int PATH '@B:bar');

```

foo	bar
1	2
3	4
4	5

(3 rows)

9.15.4. Mapping Tables to XML

The following functions map the contents of relational tables to XML values. They can be thought of as XML export functionality:

```

table_to_xml ( table regclass, nulls boolean,
               tableforest boolean, targetns text ) → xml
query_to_xml ( query text, nulls boolean,
               tableforest boolean, targetns text ) → xml
cursor_to_xml ( cursor refcursor, count integer, nulls boolean,
               tableforest boolean, targetns text ) → xml

```

`table_to_xml` maps the content of the named table, passed as parameter *table*. The *regclass* type accepts strings identifying tables using the usual notation, including optional schema qualification and double quotes (see [Section 8.19](#) for details). `query_to_xml` executes the query whose text is passed as parameter *query* and maps the result set. `cursor_to_xml` fetches the indicated number of rows from the cursor specified by the parameter *cursor*. This variant is recommended if large tables have to be mapped, because the result value is built up in memory by each function.

If *tableforest* is false, then the resulting XML document looks like this:

```

<tablename>
  <row>
    <columnname1>data</columnname1>
    <columnname2>data</columnname2>
  </row>

```

```
<row>
  ...
</row>
```

```
...
</tablename>
```

If *tableforest* is true, the result is an XML content fragment that looks like this:

```
<tablename>
  <columnname1>data</columnname1>
  <columnname2>data</columnname2>
</tablename>
```

```
<tablename>
  ...
</tablename>
```

```
...
```

If no table name is available, that is, when mapping a query or a cursor, the string *table* is used in the first format, *row* in the second format.

The choice between these formats is up to the user. The first format is a proper XML document, which will be important in many applications. The second format tends to be more useful in the *cursor_to_xml* function if the result values are to be reassembled into one document later on. The functions for producing XML content discussed above, in particular *xmlelement*, can be used to alter the results to taste.

The data values are mapped in the same way as described for the function *xmlelement* above.

The parameter *nulls* determines whether null values should be included in the output. If true, null values in columns are represented as:

```
<columnname xsi:nil="true"/>
```

where *xsi* is the XML namespace prefix for XML Schema Instance. An appropriate namespace declaration will be added to the result value. If false, columns containing null values are simply omitted from the output.

The parameter *targetns* specifies the desired XML namespace of the result. If no particular namespace is wanted, an empty string should be passed.

The following functions return XML Schema documents describing the mappings performed by the corresponding functions above:

```
table_to_xmlschema ( table regclass, nulls boolean,
                    tableforest boolean, targetns text ) → xml
query_to_xmlschema ( query text, nulls boolean,
                    tableforest boolean, targetns text ) → xml
cursor_to_xmlschema ( cursor refcursor, nulls boolean,
                    tableforest boolean, targetns text ) → xml
```

It is essential that the same parameters are passed in order to obtain matching XML data mappings and XML Schema documents.

The following functions produce XML data mappings and the corresponding XML Schema in one document (or forest), linked together. They can be useful where self-contained and self-describing results are wanted:

```
table_to_xml_and_xmlschema ( table regclass, nulls boolean,
                             tableforest boolean, targetns text ) → xml
```

```
query_to_xml_and_xmlschema ( query text, nulls boolean,
                             tableforest boolean, targetns text ) → xml
```

In addition, the following functions are available to produce analogous mappings of entire schemas or the entire current database:

```
schema_to_xml ( schema name, nulls boolean,
                tableforest boolean, targetns text ) → xml
schema_to_xmlschema ( schema name, nulls boolean,
                     tableforest boolean, targetns text ) → xml
schema_to_xml_and_xmlschema ( schema name, nulls boolean,
                             tableforest boolean, targetns text ) → xml

database_to_xml ( nulls boolean,
                 tableforest boolean, targetns text ) → xml
database_to_xmlschema ( nulls boolean,
                       tableforest boolean, targetns text ) → xml
database_to_xml_and_xmlschema ( nulls boolean,
                               tableforest boolean, targetns text ) → xml
```

These functions ignore tables that are not readable by the current user. The database-wide functions additionally ignore schemas that the current user does not have `USAGE` (lookup) privilege for.

Note that these potentially produce a lot of data, which needs to be built up in memory. When requesting content mappings of large schemas or databases, it might be worthwhile to consider mapping the tables separately instead, possibly even through a cursor.

The result of a schema content mapping looks like this:

```
<schemaname>
table1-mapping
table2-mapping
...
</schemaname>
```

where the format of a table mapping depends on the `tableforest` parameter as explained above.

The result of a database content mapping looks like this:

```
<dbname>

<schema1name>
...
</schema1name>

<schema2name>
...
</schema2name>

...

</dbname>
```

where the schema mapping is as above.

As an example of using the output produced by these functions, [Example 9.1](#) shows an XSLT stylesheet that converts the output of `table_to_xml_and_xmlschema` to an HTML document containing a tabular

rendition of the table data. In a similar manner, the results from these functions can be converted into other XML-based formats.

Example 9.1. XSLT Stylesheet for Converting SQL/XML Output to HTML

```
<?xml version="1.0"?>
<xsl:stylesheet version="1.0"
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema"
  xmlns="http://www.w3.org/1999/xhtml"
>

  <xsl:output method="xml"
    doctype-system="http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd"
    doctype-public="-//W3C/DTD XHTML 1.0 Strict//EN"
    indent="yes"/>

  <xsl:template match="/*">
    <xsl:variable name="schema" select="//xsd:schema"/>
    <xsl:variable name="tabletypename"
      select="$schema/xsd:element[@name=name(current())]/@type"/>
    <xsl:variable name="rowtypename"
      select="$schema/xsd:complexType[@name=$tabletypename]/xsd:sequence/
xsd:element[@name='row']/@type"/>

    <html>
      <head>
        <title><xsl:value-of select="name(current())"/></title>
      </head>
      <body>
        <table>
          <tr>
            <xsl:for-each select="$schema/xsd:complexType[@name=$rowtypename]/
xsd:sequence/xsd:element/@name">
              <th><xsl:value-of select="."/></th>
            </xsl:for-each>
          </tr>

          <xsl:for-each select="row">
            <tr>
              <xsl:for-each select="*">
                <td><xsl:value-of select="."/></td>
              </xsl:for-each>
            </tr>
          </xsl:for-each>
        </table>
      </body>
    </html>
  </xsl:template>

</xsl:stylesheet>
```

9.16. JSON Functions and Operators

This section describes:

- functions and operators for processing and creating JSON data
- the SQL/JSON path language
- the SQL/JSON query functions

To provide native support for JSON data types within the SQL environment, Postgres Pro implements the *SQL/JSON data model*. This model comprises sequences of items. Each item can hold SQL scalar values, with an additional SQL/JSON null value, and composite data structures that use JSON arrays and objects. The model is a formalization of the implied data model in the JSON specification [RFC 7159](#).

SQL/JSON allows you to handle JSON data alongside regular SQL data, with transaction support, including:

- Uploading JSON data into the database and storing it in regular SQL columns as character or binary strings.
- Generating JSON objects and arrays from relational data.
- Querying JSON data using SQL/JSON query functions and SQL/JSON path language expressions.

To learn more about the SQL/JSON standard, see [sqltr-19075-6](#). For details on JSON types supported in Postgres Pro, see [Section 8.14](#).

9.16.1. Processing and Creating JSON Data

Note

Functions manipulating JSONB do not accept the `'\u0000'` character. To handle this, you can specify a unicode character in the [unicode_nul_character_replacement_in_jsonb](#) configuration parameter to replace this character on the fly.

[Table 9.45](#) shows the operators that are available for use with JSON data types (see [Section 8.14](#)). In addition, the usual comparison operators shown in [Table 9.1](#) are available for `jsonb`, though not for `json`. The comparison operators follow the ordering rules for B-tree operations outlined in [Section 8.14.4](#). See also [Section 9.21](#) for the aggregate function `json_agg` which aggregates record values as JSON, the aggregate function `json_object_agg` which aggregates pairs of values into a JSON object, and their `jsonb` equivalents, `jsonb_agg` and `jsonb_object_agg`.

Table 9.45. json and jsonb Operators

Operator	Description	Example(s)
<code>json -> integer -> json</code> <code>jsonb -> integer -> jsonb</code>	Extracts <i>n</i> 'th element of JSON array (array elements are indexed from zero, but negative integers count from the end).	<code>'[{"a": "foo"}, {"b": "bar"}, {"c": "baz"}]'::json -> 2 -> {"c": "baz"}</code> <code>'[{"a": "foo"}, {"b": "bar"}, {"c": "baz"}]'::json -> -3 -> {"a": "foo"}</code>
<code>json -> text -> json</code> <code>jsonb -> text -> jsonb</code>	Extracts JSON object field with the given key.	<code>'{"a": {"b": "foo"}}'::json -> 'a' -> {"b": "foo"}</code>
<code>json ->> integer -> text</code> <code>jsonb ->> integer -> text</code>	Extracts <i>n</i> 'th element of JSON array, as text.	<code>'[1,2,3]'::json ->> 2 -> 3</code>
<code>json ->> text -> text</code> <code>jsonb ->> text -> text</code>	Extracts JSON object field with the given key, as text.	

Operator	Description	Example(s)
		<code>'{"a":1, "b":2}':::json ->> 'b' → 2</code>
		<code>json #> text[] → json</code> <code>jsonb #> text[] → jsonb</code> Extracts JSON sub-object at the specified path, where path elements can be either field keys or array indexes. <code>'{"a": { "b": ["foo", "bar"] } }':::json #> '{a,b,1}' → "bar"</code>
		<code>json #>> text[] → text</code> <code>jsonb #>> text[] → text</code> Extracts JSON sub-object at the specified path as <code>text</code> . <code>'{"a": { "b": ["foo", "bar"] } }':::json #>> '{a,b,1}' → bar</code>

Note

The field/element/path extraction operators return NULL, rather than failing, if the JSON input does not have the right structure to match the request; for example if no such key or array element exists.

Some further operators exist only for `jsonb`, as shown in [Table 9.46](#). [Section 8.14.4](#) describes how these operators can be used to effectively search indexed `jsonb` data.

Table 9.46. Additional `jsonb` Operators

Operator	Description	Example(s)
		<code>jsonb @> jsonb → boolean</code> Does the first JSON value contain the second? (See Section 8.14.3 for details about containment.) <code>'{"a":1, "b":2}':::jsonb @> '{"b":2}':::jsonb → t</code>
		<code>jsonb <@ jsonb → boolean</code> Is the first JSON value contained in the second? <code>'{"b":2}':::jsonb <@ '{"a":1, "b":2}':::jsonb → t</code>
		<code>jsonb ? text → boolean</code> Does the text string exist as a top-level key or array element within the JSON value? <code>'{"a":1, "b":2}':::jsonb ? 'b' → t</code> <code>'["a", "b", "c"]':::jsonb ? 'b' → t</code>
		<code>jsonb ? text[] → boolean</code> Do any of the strings in the text array exist as top-level keys or array elements? <code>'{"a":1, "b":2, "c":3}':::jsonb ? array['b', 'd'] → t</code>
		<code>jsonb ?& text[] → boolean</code> Do all of the strings in the text array exist as top-level keys or array elements? <code>'["a", "b", "c"]':::jsonb ?& array['a', 'b'] → t</code>
		<code>jsonb jsonb → jsonb</code> Concatenates two <code>jsonb</code> values. Concatenating two arrays generates an array containing all the elements of each input. Concatenating two objects generates an object containing the union of their keys, taking the second object's value when there are duplicate keys. All oth-

Operator	Description
Example(s)	
	<p>er cases are treated by converting a non-array input into a single-element array, and then proceeding as for two arrays. Does not operate recursively: only the top-level array or object structure is merged.</p> <pre>'["a", "b"]'::jsonb '{"a", "d"}'::jsonb → ["a", "b", "a", "d"] '{"a": "b"}'::jsonb '{"c": "d"}'::jsonb → {"a": "b", "c": "d"} '[1, 2]'::jsonb '3'::jsonb → [1, 2, 3] '{"a": "b"}'::jsonb '42'::jsonb → [{"a": "b"}, 42]</pre> <p>To append an array to another array as a single entry, wrap it in an additional layer of array, for example:</p> <pre>'[1, 2]'::jsonb jsonb_build_array('[3, 4]'::jsonb) → [1, 2, [3, 4]]</pre>
<code>jsonb - text → jsonb</code>	<p>Deletes a key (and its value) from a JSON object, or matching string value(s) from a JSON array.</p> <pre>'{"a": "b", "c": "d"}'::jsonb - 'a' → {"c": "d"} '["a", "b", "c", "b"]'::jsonb - 'b' → ["a", "c"]</pre>
<code>jsonb - text[] → jsonb</code>	<p>Deletes all matching keys or array elements from the left operand.</p> <pre>'{"a": "b", "c": "d"}'::jsonb - '{a,c}'::text[] → {}</pre>
<code>jsonb - integer → jsonb</code>	<p>Deletes the array element with specified index (negative integers count from the end). Throws an error if JSON value is not an array.</p> <pre>'["a", "b"]'::jsonb - 1 → ["a"]</pre>
<code>jsonb #- text[] → jsonb</code>	<p>Deletes the field or array element at the specified path, where path elements can be either field keys or array indexes.</p> <pre>'["a", {"b":1}]'::jsonb #- '{1,b}' → ["a", {}]</pre>
<code>jsonb @? jsonpath → boolean</code>	<p>Does JSON path return any item for the specified JSON value? (This is useful only with SQL-standard JSON path expressions, not predicate check expressions, since those always return a value.)</p> <pre>'{"a":[1,2,3,4,5]}'::jsonb @? '\$.a[*] ? (@ > 2)' → t</pre>
<code>jsonb @@ jsonpath → boolean</code>	<p>Returns the result of a JSON path predicate check for the specified JSON value. (This is useful only with predicate check expressions, not SQL-standard JSON path expressions, since it will return NULL if the path result is not a single boolean value.)</p> <pre>'{"a":[1,2,3,4,5]}'::jsonb @@ '\$.a[*] > 2' → t</pre>

Note

The `jsonpath` operators `@?` and `@@` suppress the following errors: missing object field or array element, unexpected JSON item type, datetime and numeric errors. The `jsonpath`-related functions described below can also be told to suppress these types of errors. This behavior might be helpful when searching JSON document collections of varying structure.

[Table 9.47](#) shows the functions that are available for constructing `json` and `jsonb` values. Some functions in this table have a `RETURNING` clause, which specifies the data type returned. It must be one of `json`,

`jsonb`, `bytea`, a character string type (`text`, `char`, or `varchar`), or a type that can be cast to `json`. By default, the `json` type is returned.

Table 9.47. JSON Creation Functions

Function	Description	Example(s)
<code>to_json (anyelement) → json</code> <code>to_jsonb (anyelement) → jsonb</code>	Converts any SQL value to <code>json</code> or <code>jsonb</code> . Arrays and composites are converted recursively to arrays and objects (multidimensional arrays become arrays of arrays in JSON). Otherwise, if there is a cast from the SQL data type to <code>json</code> , the cast function will be used to perform the conversion; otherwise, a scalar JSON value is produced. For any scalar other than a number, a Boolean, or a null value, the text representation will be used, with escaping as necessary to make it a valid JSON string value.	<code>to_json('Fred said "Hi."':text) → "Fred said \"Hi.\""</code> <code>to_jsonb(row(42, 'Fred said "Hi."':text)) → {"f1": 42, "f2": "Fred said \"Hi.\""}"</code>
<code>array_to_json (anyarray [, boolean]) → json</code>	Converts an SQL array to a JSON array. The behavior is the same as <code>to_json</code> except that line feeds will be added between top-level array elements if the optional boolean parameter is true.	<code>array_to_json('{{1,5},{99,100}}':int[]) → [[1,5],[99,100]]</code>
<code>json_array ([{ value_expression [FORMAT JSON] } [, ...]] [{ NULL ABSENT } ON NULL] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</code> <code>json_array ([query_expression] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</code>	Constructs a JSON array from either a series of <code>value_expression</code> parameters or from the results of <code>query_expression</code> , which must be a SELECT query returning a single column. If ABSENT ON NULL is specified, NULL values are ignored. This is always the case if a <code>query_expression</code> is used.	<code>json_array(1,true,json '{"a":null}')</code> → <code>[1, true, {"a":null}]</code> <code>json_array(SELECT * FROM (VALUES(1),(2)) t)</code> → <code>[1, 2]</code>
<code>row_to_json (record [, boolean]) → json</code>	Converts an SQL composite value to a JSON object. The behavior is the same as <code>to_json</code> except that line feeds will be added between top-level elements if the optional boolean parameter is true.	<code>row_to_json(row(1,'foo')) → {"f1":1,"f2":"foo"}</code>
<code>json_build_array (VARIADIC "any") → json</code> <code>jsonb_build_array (VARIADIC "any") → jsonb</code>	Builds a possibly-heterogeneously-typed JSON array out of a variadic argument list. Each argument is converted as per <code>to_json</code> or <code>to_jsonb</code> .	<code>json_build_array(1, 2, 'foo', 4, 5) → [1, 2, "foo", 4, 5]</code>
<code>json_build_object (VARIADIC "any") → json</code> <code>jsonb_build_object (VARIADIC "any") → jsonb</code>	Builds a JSON object out of a variadic argument list. By convention, the argument list consists of alternating keys and values. Key arguments are coerced to text; value arguments are converted as per <code>to_json</code> or <code>to_jsonb</code> .	<code>json_build_object('foo', 1, 2, row(3,'bar')) → {"foo" : 1, "2" : {"f1":3,"f2":"bar"}}</code>

Function	Description	Example(s)
<code>json_object</code>	Constructs a JSON object of all the key/value pairs given, or an empty object if none are given. <i>key_expression</i> is a scalar expression defining the JSON key, which is converted to the text type. It cannot be NULL nor can it belong to a type that has a cast to the json type. If WITH UNIQUE KEYS is specified, there must not be any duplicate <i>key_expression</i> . Any pair for which the <i>value_expression</i> evaluates to NULL is omitted from the output if ABSENT ON NULL is specified; if NULL ON NULL is specified or the clause omitted, the key is included with value NULL.	<pre>json_object ([{ key_expression { VALUE ':' } value_expression [FORMAT JSON [ENCODING UTF8]] }, ...] [{ NULL ABSENT } ON NULL] [{ WITH WITHOUT } UNIQUE [KEYS]] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</pre> <pre>json_object('code' VALUE 'P123', 'title': 'Jaws') → {"code" : "P123", "title" : "Jaws"}</pre>
<code>json_object</code> <code>jsonb_object</code>	Builds a JSON object out of a text array. The array must have either exactly one dimension with an even number of members, in which case they are taken as alternating key/value pairs, or two dimensions such that each inner array has exactly two elements, which are taken as a key/value pair. All values are converted to JSON strings.	<pre>json_object (text[]) → json</pre> <pre>jsonb_object (text[]) → jsonb</pre> <pre>json_object ('{a, 1, b, "def", c, 3.5}')</pre> <pre>→ {"a" : "1", "b" : "def", "c" : "3.5"}</pre> <pre>json_object ('{{a, 1}, {b, "def"}, {c, 3.5}}')</pre> <pre>→ {"a" : "1", "b" : "def", "c" : "3.5"}</pre>
<code>json_object</code> <code>jsonb_object</code>	This form of <code>json_object</code> takes keys and values pairwise from separate text arrays. Otherwise it is identical to the one-argument form.	<pre>json_object (keys text[], values text[]) → json</pre> <pre>jsonb_object (keys text[], values text[]) → jsonb</pre> <pre>json_object ('{a,b}', '{1,2}')</pre> <pre>→ {"a": "1", "b": "2"}</pre>
<code>json</code>	Converts a given expression specified as text or bytea string (in UTF8 encoding) into a JSON value. If <i>expression</i> is NULL, an SQL null value is returned. If WITH UNIQUE is specified, the <i>expression</i> must not contain any duplicate object keys.	<pre>json (expression [FORMAT JSON [ENCODING UTF8]] [{ WITH WITHOUT } UNIQUE [KEYS]]) → json</pre> <pre>json ('{"a":123, "b":[true,"foo"], "a":"bar"}')</pre> <pre>→ {"a":123, "b":[true, "foo"], "a":"bar"}</pre>
<code>json_scalar</code>	Converts a given SQL scalar value into a JSON scalar value. If the input is NULL, an SQL null is returned. If the input is number or a boolean value, a corresponding JSON number or boolean value is returned. For any other value, a JSON string is returned.	<pre>json_scalar (expression)</pre> <pre>json_scalar(123.45) → 123.45</pre> <pre>json_scalar(CURRENT_TIMESTAMP) → "2022-05-10T10:51:04.62128-04:00"</pre>
<code>json_serialize</code>	Converts an SQL/JSON expression into a character or binary string. The <i>expression</i> can be of any JSON type, any character string type, or bytea in UTF8 encoding. The returned type used in RETURNING can be any character string type or bytea. The default is text.	<pre>json_serialize (expression [FORMAT JSON [ENCODING UTF8]] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</pre> <pre>json_serialize ('{ "a" : 1 } ' RETURNING bytea) → \x7b20226122203a2031207d20</pre>

Table 9.48 details SQL/JSON facilities for testing JSON.

Table 9.48. SQL/JSON Testing Functions

Function signature

Description

Example(s)

```
expression IS [ NOT ] JSON [ { VALUE | SCALAR | ARRAY | OBJECT } ] [ { WITH | WITHOUT } UNIQUE [ KEYS ] ]
```

This predicate tests whether *expression* can be parsed as JSON, possibly of a specified type. If SCALAR or ARRAY or OBJECT is specified, the test is whether or not the JSON is of that particular type. If WITH UNIQUE KEYS is specified, then any object in the *expression* is also tested to see if it has duplicate keys.

```
SELECT js,
       js IS JSON "json?",
       js IS JSON SCALAR "scalar?",
       js IS JSON OBJECT "object?",
       js IS JSON ARRAY "array?"
FROM (VALUES
      ('123'), ('"abc"'), ('{"a": "b"}'), ('[1,2]'), ('abc')) foo(js);
```

js	json?	scalar?	object?	array?
123	t	t	f	f
"abc"	t	t	f	f
{"a": "b"}	t	f	t	f
[1,2]	t	f	f	t
abc	f	f	f	f

```
SELECT js,
       js IS JSON OBJECT "object?",
       js IS JSON ARRAY "array?",
       js IS JSON ARRAY WITH UNIQUE KEYS "array w. UK?",
       js IS JSON ARRAY WITHOUT UNIQUE KEYS "array w/o UK?"
FROM (VALUES ('[{"a": "1"}, {"b": "2"}, {"b": "3"}]')) foo(js);
```

```
-[ RECORD 1 ]-+-----
js           | [{"a": "1"}, {"b": "2"}, {"b": "3"}]
object?      | f
array?       | t
array w. UK? | f
array w/o UK?| t
```

Table 9.49 shows the functions that are available for processing json and jsonb values.

Table 9.49. JSON Processing Functions

Function
Description
Example(s)
<pre> json_array_elements (json) → setof json jsonb_array_elements (jsonb) → setof jsonb </pre> <p>Expands the top-level JSON array into a set of JSON values.</p> <pre> select * from json_array_elements('[1,true, [2,false]]') value ----- 1 true [2,false] </pre>

Function	Description	Example(s)
<code>json_array_elements_text (json) → setof text</code> <code>jsonb_array_elements_text (jsonb) → setof text</code>	Expands the top-level JSON array into a set of text values.	<pre>select * from json_array_elements_text('["foo", "bar"]') → value ----- foo bar</pre>
<code>json_array_length (json) → integer</code> <code>jsonb_array_length (jsonb) → integer</code>	Returns the number of elements in the top-level JSON array.	<pre>json_array_length('[1,2,3,{"f1":1,"f2":[5,6]},4]') → 5 jsonb_array_length('[]') → 0</pre>
<code>json_each (json) → setof record (key text, value json)</code> <code>jsonb_each (jsonb) → setof record (key text, value jsonb)</code>	Expands the top-level JSON object into a set of key/value pairs.	<pre>select * from json_each('{ "a": "foo", "b": "bar" }') → key value -----+----- a "foo" b "bar"</pre>
<code>json_each_text (json) → setof record (key text, value text)</code> <code>jsonb_each_text (jsonb) → setof record (key text, value text)</code>	Expands the top-level JSON object into a set of key/value pairs. The returned values will be of type text.	<pre>select * from json_each_text('{ "a": "foo", "b": "bar" }') → key value -----+----- a foo b bar</pre>
<code>json_extract_path (from_json json, VARIADIC path_elems text[]) → json</code> <code>jsonb_extract_path (from_json jsonb, VARIADIC path_elems text[]) → jsonb</code>	Extracts JSON sub-object at the specified path. (This is functionally equivalent to the #> operator, but writing the path out as a variadic list can be more convenient in some cases.)	<pre>json_extract_path('{ "f2": {"f3": 1}, "f4": {"f5": 99, "f6": "foo"} }', 'f4', 'f6') → "foo"</pre>
<code>json_extract_path_text (from_json json, VARIADIC path_elems text[]) → text</code> <code>jsonb_extract_path_text (from_json jsonb, VARIADIC path_elems text[]) → text</code>	Extracts JSON sub-object at the specified path as text. (This is functionally equivalent to the #>> operator.)	<pre>json_extract_path_text('{ "f2": {"f3": 1}, "f4": {"f5": 99, "f6": "foo"} }', 'f4', 'f6') → foo</pre>
<code>json_object_keys (json) → setof text</code> <code>jsonb_object_keys (jsonb) → setof text</code>	Returns the set of keys in the top-level JSON object.	

Function	Description	Example(s)
		<pre>select * from json_object_keys ('{"f1": "abc", "f2": {"f3": "a", "f4": "b"}}') json_object_keys ----- f1 f2</pre>
<pre>json_populate_record (base anyelement, from_json json) → anyelement jsonb_populate_record (base anyelement, from_json jsonb) → anyelement</pre>	<p>Expands the top-level JSON object to a row having the composite type of the <i>base</i> argument. The JSON object is scanned for fields whose names match column names of the output row type, and their values are inserted into those columns of the output. (Fields that do not correspond to any output column name are ignored.) In typical use, the value of <i>base</i> is just <code>NULL</code>, which means that any output columns that do not match any object field will be filled with nulls. However, if <i>base</i> isn't <code>NULL</code> then the values it contains will be used for unmatched columns.</p> <p>To convert a JSON value to the SQL type of an output column, the following rules are applied in sequence:</p> <ul style="list-style-type: none"> • A JSON null value is converted to an SQL null in all cases. • If the output column is of type <code>json</code> or <code>jsonb</code>, the JSON value is just reproduced exactly. • If the output column is a composite (row) type, and the JSON value is a JSON object, the fields of the object are converted to columns of the output row type by recursive application of these rules. • Likewise, if the output column is an array type and the JSON value is a JSON array, the elements of the JSON array are converted to elements of the output array by recursive application of these rules. • Otherwise, if the JSON value is a string, the contents of the string are fed to the input conversion function for the column's data type. • Otherwise, the ordinary text representation of the JSON value is fed to the input conversion function for the column's data type. <p>While the example below uses a constant JSON value, typical use would be to reference a <code>json</code> or <code>jsonb</code> column laterally from another table in the query's <code>FROM</code> clause. Writing <code>json_populate_record</code> in the <code>FROM</code> clause is good practice, since all of the extracted columns are available for use without duplicate function calls.</p>	<pre>create type subrowtype as (d int, e text); create type myrowtype as (a int, b text[], c subrowtype); select * from json_populate_record(null::myrowtype, '{"a": 1, "b": ["2", "a b"], "c": {"d": 4, "e": "a b c"}, "x": "foo"}') → a b c ---+-----+----- 1 {2, "a b"} (4, "a b c")</pre>
<pre>jsonb_populate_record_valid (base anyelement, from_json json) → boolean</pre>	<p>Function for testing <code>jsonb_populate_record</code>. Returns true if the input <code>jsonb_populate_record</code> would finish without an error for the given input JSON object; that is, it's valid input, false otherwise.</p>	<pre>create type jsb_char2 as (a char(2)); select jsonb_populate_record_valid(NULL::jsb_char2, '{"a": "aaa"}'); → jsonb_populate_record_valid ----- f (1 row)</pre>

Function	Description	Example(s)
		<pre>select * from jsonb_populate_record(NULL::jsb_char2, '{"a": "aaa"}') q; → ERROR: value too long for type character(2) select jsonb_populate_record_valid(NULL::jsb_char2, '{"a": "aa"}'); → jsonb_populate_record_valid ----- t (1 row) select * from jsonb_populate_record(NULL::jsb_char2, '{"a": "aa"}') q; → a ----- aa (1 row)</pre>
<pre>json_populate_recordset (base anyelement, from_json json) → setof anyelement jsonb_populate_recordset (base anyelement, from_json jsonb) → setof anyelement</pre>	<p>Expands the top-level JSON array of objects to a set of rows having the composite type of the <i>base</i> argument. Each element of the JSON array is processed as described above for <code>json[b]_populate_record</code> .</p>	<pre>create type twoints as (a int, b int); select * from json_populate_recordset(null::twoints, ' [{ "a":1, "b":2}, {"a":3, "b":4}] ') → a b ---+--- 1 2 3 4</pre>
<pre>json_to_record (json) → record jsonb_to_record (jsonb) → record</pre>	<p>Expands the top-level JSON object to a row having the composite type defined by an AS clause. (As with all functions returning <code>record</code>, the calling query must explicitly define the structure of the record with an AS clause.) The output record is filled from fields of the JSON object, in the same way as described above for <code>json[b]_populate_record</code> . Since there is no input record value, unmatched columns are always filled with nulls.</p>	<pre>create type myrowtype as (a int, b text); select * from json_to_record('{"a":1, "b":[1,2,3], "c":[1,2,3], "e":"bar", "r": { "a": 123, "b": "a b c" }}') as x(a int, b text, c int[], d text, r myrowtype) → a b c d r ---+-----+-----+---+----- 1 [1,2,3] {1,2,3} (123, "a b c")</pre>
<pre>json_to_recordset (json) → setof record jsonb_to_recordset (jsonb) → setof record</pre>	<p>Expands the top-level JSON array of objects to a set of rows having the composite type defined by an AS clause. (As with all functions returning <code>record</code>, the calling query must explicitly define the structure of the record with an AS clause.) Each element of the JSON array is processed as described above for <code>json[b]_populate_record</code> .</p>	<pre>select * from json_to_recordset(' [{ "a":1, "b":"foo"}, {"a":"2", "b":"bar"}] ') as x(a int, b text) →</pre>

Function	Description	Example(s)
		<pre> a b ---+----- 1 foo 2 </pre>
<code>jsonb_set (target jsonb, path text[], new_value jsonb[, create_if_missing boolean]) → jsonb</code>	Returns <i>target</i> with the item designated by <i>path</i> replaced by <i>new_value</i> , or with <i>new_value</i> added if <i>create_if_missing</i> is true (which is the default) and the item designated by <i>path</i> does not exist. All earlier steps in the path must exist, or the <i>target</i> is returned unchanged. As with the path oriented operators, negative integers that appear in the <i>path</i> count from the end of JSON arrays. If the last path step is an array index that is out of range, and <i>create_if_missing</i> is true, the new value is added at the beginning of the array if the index is negative, or at the end of the array if it is positive.	<pre> jsonb_set(' [{"f1":1,"f2":null},2,null,3]', '{0,f1}', '[2,3,4]', false) → [{"f1": [2, 3, 4], "f2": null}, 2, null, 3] jsonb_set(' [{"f1":1,"f2":null},2]', '{0,f3}', '[2,3,4]') → [{"f1": 1, "f2": null, "f3": [2, 3, 4]}, 2] </pre>
<code>jsonb_set_lax (target jsonb, path text[], new_value jsonb[, create_if_missing boolean [, null_value_treatment text]]) → jsonb</code>	If <i>new_value</i> is not NULL, behaves identically to <code>jsonb_set</code> . Otherwise behaves according to the value of <i>null_value_treatment</i> which must be one of 'raise_exception' , 'use_json_null' , 'delete_key' , or 'return_target' . The default is 'use_json_null' .	<pre> jsonb_set_lax(' [{"f1":1,"f2":null},2,null,3]', '{0,f1}', null) → [{"f1": null, "f2": null}, 2, null, 3] jsonb_set_lax(' [{"f1":99,"f2":null},2]', '{0,f3}', null, true, 'return_target') → [{"f1": 99, "f2": null}, 2] </pre>
<code>jsonb_insert (target jsonb, path text[], new_value jsonb[, insert_after boolean]) → jsonb</code>	Returns <i>target</i> with <i>new_value</i> inserted. If the item designated by the <i>path</i> is an array element, <i>new_value</i> will be inserted before that item if <i>insert_after</i> is false (which is the default), or after it if <i>insert_after</i> is true. If the item designated by the <i>path</i> is an object field, <i>new_value</i> will be inserted only if the object does not already contain that key. All earlier steps in the path must exist, or the <i>target</i> is returned unchanged. As with the path oriented operators, negative integers that appear in the <i>path</i> count from the end of JSON arrays. If the last path step is an array index that is out of range, the new value is added at the beginning of the array if the index is negative, or at the end of the array if it is positive.	<pre> jsonb_insert('{"a": [0,1,2]}', '{a, 1}', '"new_value"') → {"a": [0, "new_value", 1, 2]} jsonb_insert('{"a": [0,1,2]}', '{a, 1}', '"new_value"', true) → {"a": [0, 1, "new_value", 2]} </pre>
<code>json_strip_nulls (json) → json</code> <code>jsonb_strip_nulls (jsonb) → jsonb</code>	Deletes all object fields that have null values from the given JSON value, recursively. Null values that are not object fields are untouched.	<pre> json_strip_nulls(' [{"f1":1, "f2":null}, 2, null, 3]') → [{"f1":1},2, null,3] </pre>
<code>jsonb_path_exists (target jsonb, path jsonpath[, vars jsonb[, silent boolean]]) → boolean</code>		

Function	Description	Example(s)
	Checks whether the JSON path returns any item for the specified JSON value. (This is useful only with SQL-standard JSON path expressions, not predicate check expressions , since those always return a value.) If the <i>vars</i> argument is specified, it must be a JSON object, and its fields provide named values to be substituted into the <i>jsonpath</i> expression. If the <i>silent</i> argument is specified and is <code>true</code> , the function suppresses the same errors as the <code>@?</code> and <code>@@</code> operators do.	<pre>jsonb_path_exists('{\"a\":[1,2,3,4,5]}', '\$.a[*] ? (@ >= \$min && @ <= \$max)', '{\"min\":2, \"max\":4}') → t</pre>
<code>jsonb_path_match</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>boolean</code> Returns the SQL boolean result of a JSON path predicate check for the specified JSON value. (This is useful only with predicate check expressions , not SQL-standard JSON path expressions, since it will either fail or return <code>NULL</code> if the path result is not a single boolean value.) The optional <i>vars</i> and <i>silent</i> arguments act the same as for <code>jsonb_path_exists</code> .	<pre>jsonb_path_match('{\"a\":[1,2,3,4,5]}', 'exists(\$.a[*] ? (@ >= \$min && @ <= \$max))', '{\"min\":2, \"max\":4}') → t</pre>
<code>jsonb_path_query</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>setof jsonb</code> Returns all JSON items returned by the JSON path for the specified JSON value. For SQL-standard JSON path expressions it returns the JSON values selected from <i>target</i> . For predicate check expressions it returns the result of the predicate check: <code>true</code> , <code>false</code> , or <code>null</code> . The optional <i>vars</i> and <i>silent</i> arguments act the same as for <code>jsonb_path_exists</code> .	<pre>select * from jsonb_path_query('{\"a\":[1,2,3,4,5]}', '\$.a[*] ? (@ >= \$min && @ <= \$max)', '{\"min\":2, \"max\":4}') → jsonb_path_query ----- 2 3 4</pre>
<code>jsonb_path_query_array</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>jsonb</code> Returns all JSON items returned by the JSON path for the specified JSON value, as a JSON array. The parameters are the same as for <code>jsonb_path_query</code> .	<pre>jsonb_path_query_array('{\"a\":[1,2,3,4,5]}', '\$.a[*] ? (@ >= \$min && @ <= \$max)', '{\"min\":2, \"max\":4}') → [2, 3, 4]</pre>
<code>jsonb_path_query_first</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>jsonb</code> Returns the first JSON item returned by the JSON path for the specified JSON value, or <code>NULL</code> if there are no results. The parameters are the same as for <code>jsonb_path_query</code> .	<pre>jsonb_path_query_first('{\"a\":[1,2,3,4,5]}', '\$.a[*] ? (@ >= \$min && @ <= \$max)', '{\"min\":2, \"max\":4}') → 2</pre>
<code>jsonb_path_exists_tz</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>boolean</code>	
<code>jsonb_path_match_tz</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>boolean</code>	
<code>jsonb_path_query_tz</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>setof jsonb</code>	
<code>jsonb_path_query_array_tz</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent boolean</i>]]) → <code>jsonb</code>	

Function	Description	Example(s)
<code>jsonb_path_query_first_tz</code>	(<i>target jsonb</i> , <i>path jsonpath</i> [, <i>vars jsonb</i> [, <i>silent</i> boolean]]) → <i>jsonb</i>	These functions act like their counterparts described above without the <code>_tz</code> suffix, except that these functions support comparisons of date/time values that require timezone-aware conversions. The example below requires interpretation of the date-only value 2015-08-02 as a timestamp with time zone, so the result depends on the current TimeZone setting. Due to this dependency, these functions are marked as <i>stable</i> , which means these functions cannot be used in indexes. Their counterparts are <i>immutable</i> , and so can be used in indexes; but they will throw errors if asked to make such comparisons. <code>jsonb_path_exists_tz('["2015-08-01 12:00:00-05"]', '\$[*] ? (@.datetime() < "2015-08-02".datetime())') → t</code>
<code>jsonb_pretty</code>	(<i>jsonb</i>) → <i>text</i> Converts the given JSON value to pretty-printed, indented text.	<code>jsonb_pretty(' [{"f1":1,"f2":null}, 2]')</code> → <pre>[{ "f1": 1, "f2": null }, 2]</pre>
<code>json_typeof</code>	(<i>json</i>) → <i>text</i>	
<code>jsonb_typeof</code>	(<i>jsonb</i>) → <i>text</i> Returns the type of the top-level JSON value as a text string. Possible types are <i>object</i> , <i>array</i> , <i>string</i> , <i>number</i> , <i>boolean</i> , and <i>null</i> . (The <i>null</i> result should not be confused with an SQL NULL; see the examples.)	<code>json_typeof('-123.4')</code> → <i>number</i> <code>json_typeof('null'::json)</code> → <i>null</i> <code>json_typeof(NULL::json) IS NULL</code> → <i>t</i>

9.16.2. The SQL/JSON Path Language

SQL/JSON path expressions specify item(s) to be retrieved from a JSON value, similarly to XPath expressions used for access to XML content. In Postgres Pro, path expressions are implemented as the `jsonpath` data type and can use any elements described in [Section 8.14.7](#).

JSON query functions and operators pass the provided path expression to the *path engine* for evaluation. If the expression matches the queried JSON data, the corresponding JSON item, or set of items, is returned. If there is no match, the result will be *NULL*, *false*, or an error, depending on the function. Path expressions are written in the SQL/JSON path language and can include arithmetic expressions and functions.

A path expression consists of a sequence of elements allowed by the `jsonpath` data type. The path expression is normally evaluated from left to right, but you can use parentheses to change the order of operations. If the evaluation is successful, a sequence of JSON items is produced, and the evaluation result is returned to the JSON query function that completes the specified computation.

To refer to the JSON value being queried (the *context item*), use the `$` variable in the path expression. The first element of a path must always be `$`. It can be followed by one or more [accessor operators](#), which go down the JSON structure level by level to retrieve sub-items of the context item. Each accessor operator acts on the result(s) of the previous evaluation step, producing zero, one, or more output items from each input item.

For example, suppose you have some JSON data from a GPS tracker that you would like to parse, such as:

```
SELECT '{
  "track": {
    "segments": [
      {
        "location": [ 47.763, 13.4034 ],
        "start time": "2018-10-14 10:05:14",
        "HR": 73
      },
      {
        "location": [ 47.706, 13.2635 ],
        "start time": "2018-10-14 10:39:21",
        "HR": 135
      }
    ]
  }
}' AS json \gset
```

(The above example can be copied-and-pasted into psql to set things up for the following examples. Then psql will expand `: 'json'` into a suitably-quoted string constant containing the JSON value.)

To retrieve the available track segments, you need to use the `.key` accessor operator to descend through surrounding JSON objects, for example:

```
=> select jsonb_path_query(:'json', '$.track.segments');

 jsonb_path_query
-----
{"HR": 73, "location": [47.763, 13.4034], "start time": "2018-10-14 10:05:14"},
{"HR": 135, "location": [47.706, 13.2635], "start time": "2018-10-14 10:39:21"}
```

To retrieve the contents of an array, you typically use the `[*]` operator. The following example will return the location coordinates for all the available track segments:

```
=> select jsonb_path_query(:'json', '$.track.segments[*].location');

 jsonb_path_query
-----
[47.763, 13.4034]
[47.706, 13.2635]
```

Here we started with the whole JSON input value (`$`), then the `.track` accessor selected the JSON object associated with the `"track"` object key, then the `.segments` accessor selected the JSON array associated with the `"segments"` key within that object, then the `[*]` accessor selected each element of that array (producing a series of items), then the `.location` accessor selected the JSON array associated with the `"location"` key within each of those objects. In this example, each of those objects had a `"location"` key; but if any of them did not, the `.location` accessor would have simply produced no output for that input item.

To return the coordinates of the first segment only, you can specify the corresponding subscript in the `[]` accessor operator. Recall that JSON array indexes are 0-relative:

```
=> select jsonb_path_query(:'json', '$.track.segments[0].location');

 jsonb_path_query
-----
[47.763, 13.4034]
```

The result of each path evaluation step can be processed by one or more of the `jsonpath` operators and methods listed in [Section 9.16.2.3](#). Each method name must be preceded by a dot. For example, you can get the size of an array:

```
=> select jsonb_path_query(:'json', '$.track.segments.size()');
jsonb_path_query
-----
2
```

More examples of using `jsonpath` operators and methods within path expressions appear below in [Section 9.16.2.3](#).

A path can also contain *filter expressions* that work similarly to the `WHERE` clause in SQL. A filter expression begins with a question mark and provides a condition in parentheses:

```
? (condition)
```

Filter expressions must be written just after the path evaluation step to which they should apply. The result of that step is filtered to include only those items that satisfy the provided condition. SQL/JSON defines three-valued logic, so the condition can produce `true`, `false`, or `unknown`. The `unknown` value plays the same role as SQL `NULL` and can be tested for with the `is unknown` predicate. Further path evaluation steps use only those items for which the filter expression returned `true`.

The functions and operators that can be used in filter expressions are listed in [Table 9.51](#). Within a filter expression, the `@` variable denotes the value being considered (i.e., one result of the preceding path step). You can write accessor operators after `@` to retrieve component items.

For example, suppose you would like to retrieve all heart rate values higher than 130. You can achieve this as follows:

```
=> select jsonb_path_query(:'json', '$.track.segments[*].HR ? (@ > 130)');
jsonb_path_query
-----
135
```

To get the start times of segments with such values, you have to filter out irrelevant segments before selecting the start times, so the filter expression is applied to the previous step, and the path used in the condition is different:

```
=> select jsonb_path_query(:'json', '$.track.segments[*] ? (@.HR > 130).\"start time\"');
jsonb_path_query
-----
"2018-10-14 10:39:21"
```

You can use several filter expressions in sequence, if required. The following example selects start times of all segments that contain locations with relevant coordinates and high heart rate values:

```
=> select jsonb_path_query(:'json', '$.track.segments[*] ? (@.location[1] < 13.4) ?
(@.HR > 130).\"start time\"');
jsonb_path_query
-----
"2018-10-14 10:39:21"
```

Using filter expressions at different nesting levels is also allowed. The following example first filters all segments by location, and then returns high heart rate values for these segments, if available:

```
=> select jsonb_path_query(:'json', '$.track.segments[*] ? (@.location[1] < 13.4).HR ?
(@ > 130)');
jsonb_path_query
-----
135
```

You can also nest filter expressions within each other. This example returns the size of the track if it contains any segments with high heart rate values, or an empty sequence otherwise:

```
=> select jsonb_path_query(:'json', '$.track ? (exists(@.segments[*] ? (@.HR >
130))) .segments.size()');
jsonb_path_query
-----
```

9.16.2.1. Deviations from the SQL Standard

Postgres Pro's implementation of the SQL/JSON path language has the following deviations from the SQL/JSON standard.

9.16.2.1.1. Boolean Predicate Check Expressions

As an extension to the SQL standard, a Postgres Pro path expression can be a Boolean predicate, whereas the SQL standard allows predicates only within filters. While SQL-standard path expressions return the relevant element(s) of the queried JSON value, predicate check expressions return the single three-valued `jsonb` result of the predicate: `true`, `false`, or `null`. For example, we could write this SQL-standard filter expression:

```
=> select jsonb_path_query('json', '$.track.segments ?(@[*].HR > 130)');
           jsonb_path_query
-----
{"HR": 135, "location": [47.706, 13.2635], "start time": "2018-10-14 10:39:21"}
```

The similar predicate check expression simply returns `true`, indicating that a match exists:

```
=> select jsonb_path_query('json', '$.track.segments[*].HR > 130');
           jsonb_path_query
-----
true
```

Note

Predicate check expressions are required in the `@@` operator (and the `jsonb_path_match` function), and should not be used with the `@?` operator (or the `jsonb_path_exists` function).

9.16.2.1.2. Regular Expression Interpretation

There are minor differences in the interpretation of regular expression patterns used in `like_regex` filters, as described in [Section 9.16.2.4](#).

9.16.2.2. Strict and Lax Modes

When you query JSON data, the path expression may not match the actual JSON data structure. An attempt to access a non-existent member of an object or element of an array is defined as a structural error. SQL/JSON path expressions have two modes of handling structural errors:

- **lax** (default) — the path engine implicitly adapts the queried data to the specified path. Any structural errors that cannot be fixed as described below are suppressed, producing no match.
- **strict** — if a structural error occurs, an error is raised.

Lax mode facilitates matching of a JSON document and path expression when the JSON data does not conform to the expected schema. If an operand does not match the requirements of a particular operation, it can be automatically wrapped as an SQL/JSON array, or unwrapped by converting its elements into an SQL/JSON sequence before performing the operation. Also, comparison operators automatically unwrap their operands in lax mode, so you can compare SQL/JSON arrays out-of-the-box. An array of size 1 is considered equal to its sole element. Automatic unwrapping is not performed when:

- The path expression contains `type()` or `size()` methods that return the type and the number of elements in the array, respectively.
- The queried JSON data contain nested arrays. In this case, only the outermost array is unwrapped, while all the inner arrays remain unchanged. Thus, implicit unwrapping can only go one level down within each path evaluation step.

For example, when querying the GPS data listed above, you can abstract from the fact that it stores an array of segments when using lax mode:

```
=> select jsonb_path_query('json', 'lax $.track.segments.location');
```

```
jsonb_path_query
-----
[47.763, 13.4034]
[47.706, 13.2635]
```

In strict mode, the specified path must exactly match the structure of the queried JSON document, so using this path expression will cause an error:

```
=> select jsonb_path_query(:'json', 'strict $.track.segments.location');
ERROR:  jsonpath member accessor can only be applied to an object
```

To get the same result as in lax mode, you have to explicitly unwrap the `segments` array:

```
=> select jsonb_path_query(:'json', 'strict $.track.segments[*].location');
jsonb_path_query
-----
[47.763, 13.4034]
[47.706, 13.2635]
```

The unwrapping behavior of lax mode can lead to surprising results. For instance, the following query using the `.**` accessor selects every `HR` value twice:

```
=> select jsonb_path_query(:'json', 'lax $.**.HR');
jsonb_path_query
-----
73
135
73
135
```

This happens because the `.**` accessor selects both the `segments` array and each of its elements, while the `.HR` accessor automatically unwraps arrays when using lax mode. To avoid surprising results, we recommend using the `.**` accessor only in strict mode. The following query selects each `HR` value just once:

```
=> select jsonb_path_query(:'json', 'strict $.**.HR');
jsonb_path_query
-----
73
135
```

The unwrapping of arrays can also lead to unexpected results. Consider this example, which selects all the `location` arrays:

```
=> select jsonb_path_query(:'json', 'lax $.track.segments[*].location');
jsonb_path_query
-----
[47.763, 13.4034]
[47.706, 13.2635]
(2 rows)
```

As expected it returns the full arrays. But applying a filter expression causes the arrays to be unwrapped to evaluate each item, returning only the items that match the expression:

```
=> select jsonb_path_query(:'json', 'lax $.track.segments[*].location ?(@[*] > 15)');
jsonb_path_query
-----
47.763
47.706
(2 rows)
```

This despite the fact that the full arrays are selected by the path expression. Use strict mode to restore selecting the arrays:

```
=> select jsonb_path_query(:'json', 'strict $.track.segments[*].location ?(@[*] > 15)');
```

```
jsonb_path_query
-----
[47.763, 13.4034]
[47.706, 13.2635]
(2 rows)
```

9.16.2.3. SQL/JSON Path Operators and Methods

[Table 9.50](#) shows the operators and methods available in `jsonpath`. Note that while the unary operators and methods can be applied to multiple values resulting from a preceding path step, the binary operators (addition etc.) can only be applied to single values. In lax mode, methods applied to an array will be executed for each value in the array. The exceptions are `.type()` and `.size()`, which apply to the array itself.

Table 9.50. `jsonpath` Operators and Methods

Operator/Method Description Example(s)
$number + number \rightarrow number$ Addition <code>jsonb_path_query('[2]', '\$[0] + 3')</code> $\rightarrow 5$
$+ number \rightarrow number$ Unary plus (no operation); unlike addition, this can iterate over multiple values <code>jsonb_path_query_array('{ "x": [2,3,4] }', '+ \$.x')</code> $\rightarrow [2, 3, 4]$
$number - number \rightarrow number$ Subtraction <code>jsonb_path_query('[2]', '7 - \$[0]')</code> $\rightarrow 5$
$- number \rightarrow number$ Negation; unlike subtraction, this can iterate over multiple values <code>jsonb_path_query_array('{ "x": [2,3,4] }', '- \$.x')</code> $\rightarrow [-2, -3, -4]$
$number * number \rightarrow number$ Multiplication <code>jsonb_path_query('[4]', '2 * \$[0]')</code> $\rightarrow 8$
$number / number \rightarrow number$ Division <code>jsonb_path_query('[8.5]', '\$[0] / 2')</code> $\rightarrow 4.2500000000000000$
$number \% number \rightarrow number$ Modulo (remainder) <code>jsonb_path_query('[32]', '\$[0] \% 10')</code> $\rightarrow 2$
$value.type() \rightarrow string$ Type of the JSON item (see <code>json_typeof</code>) <code>jsonb_path_query_array('[1, "2", {}]', '\$[*].type()')</code> $\rightarrow ["number", "string", "object"]$
$value.size() \rightarrow number$ Size of the JSON item (number of array elements, or 1 if not an array) <code>jsonb_path_query('{ "m": [11, 15] }', '\$.m.size()')</code> $\rightarrow 2$
$value.boolean() \rightarrow boolean$ Boolean value converted from a JSON boolean, number, or string

Operator/Method Description Example(s)
<pre>jsonb_path_query_array('[1, "yes", false]', '\$[*].boolean()') → [true, true, false]</pre>
<pre>value.string() → string</pre> <p>String value converted from a JSON boolean, number, string, or datetime</p> <pre>jsonb_path_query_array('[1.23, "xyz", false]', '\$[*].string()') → ["1.23", "xyz", "false"]</pre> <pre>jsonb_path_query('"2023-08-15 12:34:56"', '\$.timestamp().string()') → "2023-08-15T12:34:56"</pre>
<pre>value.double() → number</pre> <p>Approximate floating-point number converted from a JSON number or string</p> <pre>jsonb_path_query('{"len": "1.9"}', '\$.len.double() * 2') → 3.8</pre>
<pre>number.ceiling() → number</pre> <p>Nearest integer greater than or equal to the given number</p> <pre>jsonb_path_query('{"h": 1.3}', '\$.h.ceiling()') → 2</pre>
<pre>number.floor() → number</pre> <p>Nearest integer less than or equal to the given number</p> <pre>jsonb_path_query('{"h": 1.7}', '\$.h.floor()') → 1</pre>
<pre>number.abs() → number</pre> <p>Absolute value of the given number</p> <pre>jsonb_path_query('{"z": -0.3}', '\$.z.abs()') → 0.3</pre>
<pre>value.bigint() → bigint</pre> <p>Big integer value converted from a JSON number or string</p> <pre>jsonb_path_query('{"len": "9876543219"}', '\$.len.bigint()') → 9876543219</pre>
<pre>value.decimal([precision [, scale]]) → decimal</pre> <p>Rounded decimal value converted from a JSON number or string (precision and scale must be integer values)</p> <pre>jsonb_path_query('1234.5678', '\$.decimal(6, 2)') → 1234.57</pre>
<pre>value.integer() → integer</pre> <p>Integer value converted from a JSON number or string</p> <pre>jsonb_path_query('{"len": "12345"}', '\$.len.integer()') → 12345</pre>
<pre>value.number() → numeric</pre> <p>Numeric value converted from a JSON number or string</p> <pre>jsonb_path_query('{"len": "123.45"}', '\$.len.number()') → 123.45</pre>
<pre>string.datetime() → datetime_type (see note)</pre> <p>Date/time value converted from a string</p> <pre>jsonb_path_query('["2015-8-1", "2015-08-12"]', '\$[*] ? (@.datetime() < "2015-08-2".datetime())') → "2015-8-1"</pre>
<pre>string.datetime(template) → datetime_type (see note)</pre> <p>Date/time value converted from a string using the specified to_timestamp template</p> <pre>jsonb_path_query_array('["12:30", "18:40"]', '\$[*].datetime("HH24:MI")') → ["12:30:00", "18:40:00"]</pre>
<pre>string.date() → date</pre> <p>Date value converted from a string</p>

Operator/Method Description Example(s)
<pre>jsonb_path_query('2023-08-15', '\$.date()') → "2023-08-15"</pre>
<pre>string.time() → time without time zone</pre> <p>Time without time zone value converted from a string</p> <pre>jsonb_path_query('12:34:56', '\$.time()') → "12:34:56"</pre>
<pre>string.time(precision) → time without time zone</pre> <p>Time without time zone value converted from a string, with fractional seconds adjusted to the given precision</p> <pre>jsonb_path_query('12:34:56.789', '\$.time(2)') → "12:34:56.79"</pre>
<pre>string.time_tz() → time with time zone</pre> <p>Time with time zone value converted from a string</p> <pre>jsonb_path_query('12:34:56 +05:30', '\$.time_tz()') → "12:34:56+05:30"</pre>
<pre>string.time_tz(precision) → time with time zone</pre> <p>Time with time zone value converted from a string, with fractional seconds adjusted to the given precision</p> <pre>jsonb_path_query('12:34:56.789 +05:30', '\$.time_tz(2)') → "12:34:56.79+05:30"</pre>
<pre>string.timestamp() → timestamp without time zone</pre> <p>Timestamp without time zone value converted from a string</p> <pre>jsonb_path_query('2023-08-15 12:34:56', '\$.timestamp()') → "2023-08-15T12:34:56"</pre>
<pre>string.timestamp(precision) → timestamp without time zone</pre> <p>Timestamp without time zone value converted from a string, with fractional seconds adjusted to the given precision</p> <pre>jsonb_path_query('2023-08-15 12:34:56.789', '\$.timestamp(2)') → "2023-08-15T12:34:56.79"</pre>
<pre>string.timestamp_tz() → timestamp with time zone</pre> <p>Timestamp with time zone value converted from a string</p> <pre>jsonb_path_query('2023-08-15 12:34:56 +05:30', '\$.timestamp_tz()') → "2023-08-15T12:34:56+05:30"</pre>
<pre>string.timestamp_tz(precision) → timestamp with time zone</pre> <p>Timestamp with time zone value converted from a string, with fractional seconds adjusted to the given precision</p> <pre>jsonb_path_query('2023-08-15 12:34:56.789 +05:30', '\$.timestamp_tz(2)') → "2023-08-15T12:34:56.79+05:30"</pre>
<pre>object.keyvalue() → array</pre> <p>The object's key-value pairs, represented as an array of objects containing three fields: "key", "value", and "id"; "id" is a unique identifier of the object the key-value pair belongs to</p> <pre>jsonb_path_query_array('{ "x": "20", "y": 32 }', '\$.keyvalue()') → [{ "id": 0, "key": "x", "value": "20"}, { "id": 0, "key": "y", "value": 32 }]</pre>

Note

The result type of the `datetime()` and `datetime(template)` methods can be `date`, `timetz`, `time`, `timestamptz`, or `timestamp`. Both methods determine their result type dynamically.

The `datetime()` method sequentially tries to match its input string to the ISO formats for `date`, `timetz`, `time`, `timestamptz`, and `timestamp`. It stops on the first matching format and emits the corresponding data type.

The `datetime(template)` method determines the result type according to the fields used in the provided template string.

The `datetime()` and `datetime(template)` methods use the same parsing rules as the `to_timestamp` SQL function does (see [Section 9.8](#)), with three exceptions. First, these methods don't allow unmatched template patterns. Second, only the following separators are allowed in the template string: minus sign, period, solidus (slash), comma, apostrophe, semicolon, colon and space. Third, separators in the template string must exactly match the input string.

If different date/time types need to be compared, an implicit cast is applied. A `date` value can be cast to `timestamp` or `timestamptz`, `timestamp` can be cast to `timestamptz`, and `time` to `timetz`. However, all but the first of these conversions depend on the current [TimeZone](#) setting, and thus can only be performed within timezone-aware `jsonpath` functions. Similarly, other date/time-related methods that convert strings to date/time types also do this casting, which may involve the current [TimeZone](#) setting. Therefore, these conversions can also only be performed within timezone-aware `jsonpath` functions.

[Table 9.51](#) shows the available filter expression elements.

Table 9.51. jsonpath Filter Expression Elements

Predicate/Value	Description	Example(s)
<code>value == value → boolean</code>	Equality comparison (this, and the other comparison operators, work on all JSON scalar values)	<code>jsonb_path_query_array('[1, "a", 1, 3]', '\$[*] ? (@ == 1)')</code> → [1, 1] <code>jsonb_path_query_array('[1, "a", 1, 3]', '\$[*] ? (@ == "a")')</code> → ["a"]
<code>value != value → boolean</code> <code>value <> value → boolean</code>	Non-equality comparison	<code>jsonb_path_query_array('[1, 2, 1, 3]', '\$[*] ? (@ != 1)')</code> → [2, 3] <code>jsonb_path_query_array('["a", "b", "c"]', '\$[*] ? (@ <> "b")')</code> → ["a", "c"]
<code>value < value → boolean</code>	Less-than comparison	<code>jsonb_path_query_array('[1, 2, 3]', '\$[*] ? (@ < 2)')</code> → [1]
<code>value <= value → boolean</code>	Less-than-or-equal-to comparison	<code>jsonb_path_query_array('["a", "b", "c"]', '\$[*] ? (@ <= "b")')</code> → ["a", "b"]
<code>value > value → boolean</code>	Greater-than comparison	<code>jsonb_path_query_array('[1, 2, 3]', '\$[*] ? (@ > 2)')</code> → [3]
<code>value >= value → boolean</code>	Greater-than-or-equal-to comparison	<code>jsonb_path_query_array('[1, 2, 3]', '\$[*] ? (@ >= 2)')</code> → [2, 3]

Predicate/Value Description Example(s)
<p><code>true</code> → <code>boolean</code> JSON constant true <code>jsonb_path_query('[{ "name": "John", "parent": false}, {"name": "Chris", "parent": true}]', '\$[*] ? (@.parent == true)') → {"name": "Chris", "parent": true}</code></p>
<p><code>false</code> → <code>boolean</code> JSON constant false <code>jsonb_path_query('[{ "name": "John", "parent": false}, {"name": "Chris", "parent": true}]', '\$[*] ? (@.parent == false)') → {"name": "John", "parent": false}</code></p>
<p><code>null</code> → <code>value</code> JSON constant null (note that, unlike in SQL, comparison to null works normally) <code>jsonb_path_query('[{ "name": "Mary", "job": null}, {"name": "Michael", "job": "driver"}]', '\$[*] ? (@.job == null) .name') → "Mary"</code></p>
<p><code>boolean && boolean</code> → <code>boolean</code> Boolean AND <code>jsonb_path_query('[1, 3, 7]', '\$[*] ? (@ > 1 && @ < 5)') → 3</code></p>
<p><code>boolean boolean</code> → <code>boolean</code> Boolean OR <code>jsonb_path_query('[1, 3, 7]', '\$[*] ? (@ < 1 @ > 5)') → 7</code></p>
<p><code>! boolean</code> → <code>boolean</code> Boolean NOT <code>jsonb_path_query('[1, 3, 7]', '\$[*] ? (!(@ < 5))') → 7</code></p>
<p><code>boolean is unknown</code> → <code>boolean</code> Tests whether a Boolean condition is unknown. <code>jsonb_path_query('[-1, 2, 7, "foo"]', '\$[*] ? ((@ > 0) is unknown)') → "foo"</code></p>
<p><code>string like_regex string [flag string]</code> → <code>boolean</code> Tests whether the first operand matches the regular expression given by the second operand, optionally with modifications described by a string of flag characters (see Section 9.16.2.4). <code>jsonb_path_query_array('["abc", "abd", "aBdC", "abdacb", "babc"]', '\$[*] ? (@ like_regex "^ab.*c")') → ["abc", "abdacb"]</code> <code>jsonb_path_query_array('["abc", "abd", "aBdC", "abdacb", "babc"]', '\$[*] ? (@ like_regex "^ab.*c" flag "i")') → ["abc", "aBdC", "abdacb"]</code></p>
<p><code>string starts with string</code> → <code>boolean</code> Tests whether the second operand is an initial substring of the first operand. <code>jsonb_path_query('["John Smith", "Mary Stone", "Bob Johnson"]', '\$[*] ? (@ starts with "John")') → "John Smith"</code></p>
<p><code>exists (path_expression)</code> → <code>boolean</code> Tests whether a path expression matches at least one SQL/JSON item. Returns unknown if the path expression would result in an error; the second example uses this to avoid a no-such-key error in strict mode. <code>jsonb_path_query('{ "x": [1, 2], "y": [2, 4] }', 'strict \$.* ? (exists (@ ? (@[*] > 2)))') → [2, 4]</code></p>

Predicate/Value	Description	Example(s)
		<code>jsonb_path_query_array('{ "value": 41}', 'strict \$? (exists (@.name)) .name') → []</code>

9.16.2.4. SQL/JSON Regular Expressions

SQL/JSON path expressions allow matching text to a regular expression with the `like_regex` filter. For example, the following SQL/JSON path query would case-insensitively match all strings in an array that start with an English vowel:

```
$[*] ? (@ like_regex "^[aeiou]" flag "i")
```

The optional `flag` string may include one or more of the characters `i` for case-insensitive match, `m` to allow `^` and `$` to match at newlines, `s` to allow `.` to match a newline, and `q` to quote the whole pattern (reducing the behavior to a simple substring match).

The SQL/JSON standard borrows its definition for regular expressions from the `LIKE_REGEX` operator, which in turn uses the XQuery standard. Postgres Pro does not currently support the `LIKE_REGEX` operator. Therefore, the `like_regex` filter is implemented using the POSIX regular expression engine described in [Section 9.7.3](#). This leads to various minor discrepancies from standard SQL/JSON behavior, which are cataloged in [Section 9.7.3.8](#). Note, however, that the flag-letter incompatibilities described there do not apply to SQL/JSON, as it translates the XQuery flag letters to match what the POSIX engine expects.

Keep in mind that the pattern argument of `like_regex` is a JSON path string literal, written according to the rules given in [Section 8.14.7](#). This means in particular that any backslashes you want to use in the regular expression must be doubled. For example, to match string values of the root document that contain only digits:

```
$.* ? (@ like_regex "^\\d+$")
```

9.16.3. SQL/JSON Query Functions

SQL/JSON functions `JSON_EXISTS()`, `JSON_QUERY()`, and `JSON_VALUE()` described in [Table 9.52](#) can be used to query JSON documents. Each of these functions apply a *path_expression* (an SQL/JSON path query) to a *context_item* (the document). See [Section 9.16.2](#) for more details on what the *path_expression* can contain. The *path_expression* can also reference variables, whose values are specified with their respective names in the `PASSING` clause that is supported by each function. *context_item* can be a `jsonb` value or a character string that can be successfully cast to `jsonb`.

Table 9.52. SQL/JSON Query Functions

Function signature	Description	Example(s)
<pre>JSON_EXISTS (context_item , path_expression [PASSING { value AS varname } [, ...]] [{ TRUE FALSE UNKNOWN ERROR } ON ERROR]) → boolean</pre>		
<ul style="list-style-type: none"> Returns true if the SQL/JSON <i>path_expression</i> applied to the <i>context_item</i> yields any items, false otherwise. The <code>ON ERROR</code> clause specifies the behavior if an error occurs during <i>path_expression</i> evaluation. Specifying <code>ERROR</code> will cause an error to be thrown with the appropriate message. Other options include returning boolean values <code>FALSE</code> or <code>TRUE</code> or the value <code>UNKNOWN</code> which is actually an SQL NULL. The default when no <code>ON ERROR</code> clause is specified is to return the boolean value <code>FALSE</code>. 		

Function signature
Description
Example(s)

Examples:

JSON_EXISTS(jsonb '{"key1": [1,2,3]}', 'strict \$.key1[*] ? (@ > \$x)'

PASSING 2 AS x) → t

JSON_EXISTS(jsonb '{"a": [1,2,3]}', 'lax \$.a[5]' ERROR ON ERROR) → f

JSON_EXISTS(jsonb '{"a": [1,2,3]}', 'strict \$.a[5]' ERROR ON ERROR) →

ERROR: jsonpath array subscript is out of bounds

```
JSON_QUERY (
    context_item , path_expression
    [ PASSING { value AS varname } [, ...] ]
    [ RETURNING data_type [ FORMAT JSON [ ENCODING UTF8 ] ] ]
    [ { WITHOUT | WITH { CONDITIONAL | [UNCONDITIONAL] } } [ ARRAY ] WRAPPER ]
    [ { KEEP | OMIT } QUOTES [ ON SCALAR STRING ] ]
    [ { ERROR | NULL | EMPTY { [ ARRAY ] | OBJECT } | DEFAULT expression } ON EMPTY
    ]
    [ { ERROR | NULL | EMPTY { [ ARRAY ] | OBJECT } | DEFAULT expression } ON ERROR
    ] ) → jsonb
```

- Returns the result of applying the SQL/JSON *path_expression* to the *context_item* .
- By default, the result is returned as a value of type `jsonb`, though the `RETURNING` clause can be used to return as some other type to which it can be successfully coerced.
- If the path expression may return multiple values, it might be necessary to wrap those values using the `WITH WRAPPER` clause to make it a valid JSON string, because the default behavior is to not wrap them, as if `WITHOUT WRAPPER` were specified. The `WITH WRAPPER` clause is by default taken to mean `WITH UNCONDITIONAL WRAPPER`, which means that even a single result value will be wrapped. To apply the wrapper only when multiple values are present, specify `WITH CONDITIONAL WRAPPER`. Getting multiple values in result will be treated as an error if `WITHOUT WRAPPER` is specified.
- If the result is a scalar string, by default, the returned value will be surrounded by quotes, making it a valid JSON value. It can be made explicit by specifying `KEEP QUOTES`. Conversely, quotes can be omitted by specifying `OMIT QUOTES`. To ensure that the result is a valid JSON value, `OMIT QUOTES` cannot be specified when `WITH WRAPPER` is also specified.
- The `ON EMPTY` clause specifies the behavior if evaluating *path_expression* yields an empty set. The `ON ERROR` clause specifies the behavior if an error occurs when evaluating *path_expression* , when coercing the result value to the `RETURNING` type, or when evaluating the `ON EMPTY` expression if the *path_expression* evaluation returns an empty set.
- For both `ON EMPTY` and `ON ERROR`, specifying `ERROR` will cause an error to be thrown with the appropriate message. Other options include returning an SQL NULL, an empty array (`EMPTY [ARRAY]`), an empty object (`EMPTY OBJECT`), or a user-specified expression (`DEFAULT expression`) that can be coerced to `jsonb` or the type specified in `RETURNING`. The default when `ON EMPTY` or `ON ERROR` is not specified is to return an SQL NULL value.

Examples:

JSON_QUERY(jsonb '[1,[2,3],null]', 'lax \$[*][\$off]' PASSING 1 AS off WITH
CONDITIONAL WRAPPER) → 3

JSON_QUERY(jsonb '{"a": "[1, 2]"}', 'lax \$.a' OMIT QUOTES) → [1, 2]

JSON_QUERY(jsonb '{"a": "[1, 2]"}', 'lax \$.a' RETURNING int[] OMIT QUOTES
ERROR ON ERROR) →

ERROR: malformed array literal: "[1, 2]"

Function signature
Description
Example(s)
DETAIL: Missing "]" after array dimensions.
<pre>JSON_VALUE (context_item , path_expression [PASSING { value AS varname } [, ...]] [RETURNING data_type] [{ ERROR NULL DEFAULT expression } ON EMPTY] [{ ERROR NULL DEFAULT expression } ON ERROR]) →text</pre> <ul style="list-style-type: none"> Returns the result of applying the SQL/JSON <i>path_expression</i> to the <i>context_item</i> . Only use <code>JSON_VALUE()</code> if the extracted value is expected to be a single SQL/JSON scalar item; getting multiple values will be treated as an error. If you expect that extracted value might be an object or an array, use the <code>JSON_QUERY</code> function instead. By default, the result, which must be a single scalar value, is returned as a value of type <code>text</code>, though the <code>RETURNING</code> clause can be used to return as some other type to which it can be successfully coerced. The <code>ON ERROR</code> and <code>ON EMPTY</code> clauses have similar semantics as mentioned in the description of <code>JSON_QUERY</code> , except the set of values returned in lieu of throwing an error is different. Note that scalar strings returned by <code>JSON_VALUE</code> always have their quotes removed, equivalent to specifying <code>OMIT QUOTES</code> in <code>JSON_QUERY</code> . <p>Examples:</p> <pre>JSON_VALUE(jsonb '"123.45"', '\$' RETURNING float) →123.45 JSON_VALUE(jsonb '"03:04 2015-02-01"', '\$.datetime("HH24:MI YYYY-MM-DD")' RETURNING date) →2015-02-01 JSON_VALUE(jsonb '[1,2]', 'strict \${\$off}' PASSING 1 as off) →2 JSON_VALUE(jsonb '[1,2]', 'strict \$[*]' DEFAULT 9 ON ERROR) →9</pre>

Note

The *context_item* expression is converted to `jsonb` by an implicit cast if the expression is not already of type `jsonb`. Note, however, that any parsing errors that occur during that conversion are thrown unconditionally, that is, are not handled according to the (specified or implicit) `ON ERROR` clause.

Note

`JSON_VALUE()` returns an SQL NULL if *path_expression* returns a JSON null, whereas `JSON_QUERY()` returns the JSON null as is.

9.16.4. JSON_TABLE

`JSON_TABLE` is an SQL/JSON function which queries JSON data and presents the results as a relational view, which can be accessed as a regular SQL table. You can use `JSON_TABLE` inside the `FROM` clause of a `SELECT`, `UPDATE`, or `DELETE` and as data source in a `MERGE` statement.

Taking JSON data as input, `JSON_TABLE` uses a JSON path expression to extract a part of the provided data to use as a *row pattern* for the constructed view. Each SQL/JSON value given by the row pattern serves as source for a separate row in the constructed view.

To split the row pattern into columns, `JSON_TABLE` provides the `COLUMNS` clause that defines the schema of the created view. For each column, a separate JSON path expression can be specified to be evaluated against the row pattern to get an SQL/JSON value that will become the value for the specified column in a given output row.

JSON data stored at a nested level of the row pattern can be extracted using the `NESTED PATH` clause. Each `NESTED PATH` clause can be used to generate one or more columns using the data from a nested level of the row pattern. Those columns can be specified using a `COLUMNS` clause that looks similar to the top-level `COLUMNS` clause. Rows constructed from `NESTED COLUMNS` are called *child rows* and are joined against the row constructed from the columns specified in the parent `COLUMNS` clause to get the row in the final view. Child columns themselves may contain a `NESTED PATH` specification thus allowing to extract data located at arbitrary nesting levels. Columns produced by multiple `NESTED PATHS` at the same level are considered to be *siblings* of each other and their rows after joining with the parent row are combined using `UNION`.

The rows produced by `JSON_TABLE` are laterally joined to the row that generated them, so you do not have to explicitly join the constructed view with the original table holding JSON data. Optionally, you can specify how to join the columns returned by `NESTED PATH` using the `PLAN` clause.

Each `NESTED PATH` clause can generate one or more columns. Columns produced by `NESTED PATHS` at the same level are considered to be *siblings*, while a column produced by a `NESTED PATH` is considered to be a child of the column produced by a `NESTED PATH` or row expression at a higher level. Sibling columns are always joined first. Once they are processed, the resulting rows are joined to the parent row.

The syntax is:

```
JSON_TABLE (
  context_item, path_expression [ AS json_path_name ] [ PASSING { value AS varname }
  [, ...] ]
  COLUMNS ( json_table_column [, ...] )
  [ { ERROR | EMPTY } ON ERROR ]
  [
    PLAN ( json_table_plan ) |
    PLAN DEFAULT ( { INNER | OUTER } [ , { CROSS | UNION } ]
                  | { CROSS | UNION } [ , { INNER | OUTER } ] )
  ]
)
```

where `json_table_column` is:

```
name type [ PATH json_path_specification ]
  [ { WITHOUT | WITH { CONDITIONAL | [UNCONDITIONAL] } } [ ARRAY ] WRAPPER ]
  [ { KEEP | OMIT } QUOTES [ ON SCALAR STRING ] ]
  [ { ERROR | NULL | EMPTY { [ARRAY] | OBJECT } | DEFAULT expression } ON EMPTY ]
  [ { ERROR | NULL | EMPTY { [ARRAY] | OBJECT } | DEFAULT expression } ON ERROR ]
| name type EXISTS [ PATH path_expression ]
  [ { ERROR | TRUE | FALSE | UNKNOWN } ON ERROR ]
| NESTED PATH json_path_specification [ AS path_name ]
  COLUMNS ( json_table_column [, ...] )
| name FOR ORDINALITY
```

`json_table_plan` is:

```
json_path_name [ { OUTER | INNER } json_table_plan_primary ]
| json_table_plan_primary { UNION json_table_plan_primary } [...]
| json_table_plan_primary { CROSS json_table_plan_primary } [...]
```

`json_table_plan_primary` is:

```
json_path_name | ( json_table_plan )
```

Each syntax element is described below in more detail.

```
context_item, path_expression [ AS json_path_name ] [ PASSING { value AS varname } [, ...]]
```

The *context_item* specifies the input document to query, the *path_expression* is an SQL/JSON path expression defining the query, and *json_path_name* is an optional name for the *path_expression*. The optional `PASSING` clause provides data values for the variables mentioned in the *path_expression*. The result of the input data evaluation using the aforementioned elements is called the *row pattern*, which is used as the source for row values in the constructed view.

```
COLUMNS ( json_table_column [, ...] )
```

The `COLUMNS` clause defining the schema of the constructed view. In this clause, you can specify each column to be filled with an SQL/JSON value obtained by applying a JSON path expression against the row pattern. *json_table_column* has the following variants:

```
name type [FORMAT JSON [ENCODING UTF8]] [ PATH path_expression ]
```

Inserts an SQL/JSON value obtained by applying *path_expression* against the row pattern into the view's output row after coercing it to specified *type*.

Specifying `FORMAT JSON` makes it explicit that you expect the value to be a valid json object. It only makes sense to specify `FORMAT JSON` if *type* is one of `bpchar`, `bytea`, `character varying`, `name`, `json`, `jsonb`, `text`, or a domain over these types.

Optionally, you can specify `WRAPPER`, `QUOTES` clauses to format the output and `ON EMPTY` and `ON ERROR` to handle those scenarios appropriately.

Note

This clause is internally turned into and has the same semantics as `JSON_VALUE` or `JSON_QUERY`. The latter if the specified *type* is not a scalar type or if either of `FORMAT JSON`, `WRAPPER`, or `QUOTES` clause is present.

```
name type EXISTS [ PATH path_expression ]
```

Inserts a boolean value obtained by applying *path_expression* against the row pattern into the view's output row after coercing it to specified *type*.

The value corresponds to whether applying the `PATH` expression to the row pattern yields any values.

The specified *type* should have a cast from the `boolean` type.

Optionally, you can use `ON ERROR` to specify whether to throw the error or return the specified value when an error occurs during JSON path evaluation or when coercing SQL/JSON value to the specified *type*. The default is to return a boolean value `FALSE`.

Note

This clause is internally turned into and has the same semantics as `JSON_EXISTS`.

```
NESTED [ PATH ] path_expression [ AS json_path_name ] COLUMNS ( json_table_column [, ...] )
```

Extracts SQL/JSON values from nested levels of the row pattern, generates one or more columns as defined by the `COLUMNS` subclause, and inserts the extracted SQL/JSON values into those columns. The *json_table_column* expression in the `COLUMNS` subclause uses the same syntax as in the parent `COLUMNS` clause.

The `NESTED PATH` syntax is recursive, so you can go down multiple nested levels by specifying several `NESTED PATH` subclauses within each other. It allows to unnest the hierarchy of JSON objects and arrays in a single function invocation rather than chaining several `JSON_TABLE` expressions in an SQL statement.

You can use the `PLAN` clause to define how to join the columns returned by `NESTED PATH` clauses.

`name FOR ORDINALITY`

Adds an ordinality column that provides sequential row numbering. You can have only one ordinality column per table. Row numbering is 1-based. For child rows that result from the `NESTED PATH` clauses, the parent row number is repeated.

Note

In each variant of `json_table_column` described above, if the `PATH` clause is omitted, path expression `$.name` is used, where `name` is the provided column name.

`AS json_path_name`

The optional `json_path_name` serves as an identifier of the provided `path_expression`. The name must be unique and distinct from the column names. When using the `PLAN` clause, you must specify the names for all the paths, including the row pattern. Each path name can appear in the `PLAN` clause only once.

`PLAN (json_table_plan)`

Defines how to join the data returned by `NESTED PATH` clauses to the constructed view.

To join columns with parent/child relationship, you can use:

`INNER`

Use `INNER JOIN`, so that the parent row is omitted from the output if it does not have any child rows after joining the data returned by `NESTED PATH`.

`OUTER`

Use `LEFT OUTER JOIN`, so that the parent row is always included into the output even if it does not have any child rows after joining the data returned by `NESTED PATH`, with `NULL` values inserted into the child columns if the corresponding values are missing.

This is the default option for joining columns with parent/child relationship.

To join sibling columns, you can use:

`UNION`

Generate one row for each value produced by each of the sibling columns. The columns from the other siblings are set to null.

This is the default option for joining sibling columns.

`CROSS`

Generate one row for each combination of values from the sibling columns.

`PLAN DEFAULT (OUTER | INNER [, UNION | CROSS])`

The terms can also be specified in reverse order. The `INNER` or `OUTER` option defines the joining plan for parent/child columns, while `UNION` or `CROSS` affects joins of sibling columns. This form of `PLAN` overrides the default plan for all columns at once. Even though the path names are not included in the `PLAN DEFAULT` form, to conform to the SQL/JSON standard they must be provided for all the paths if the `PLAN` clause is used.

`PLAN DEFAULT` is simpler than specifying a complete `PLAN`, and is often all that is required to get the desired output.

```
{ ERROR | EMPTY } ON ERROR
```

The optional `ON ERROR` can be used to specify how to handle errors when evaluating the top-level *path_expression*. Use `ERROR` if you want the errors to be thrown and `EMPTY` to return an empty table, that is, a table containing 0 rows. Note that this clause does not affect the errors that occur when evaluating columns, for which the behavior depends on whether the `ON ERROR` clause is specified against a given column.

Examples

In the examples that follow, the following table containing JSON data will be used:

```
CREATE TABLE my_films ( js jsonb );

INSERT INTO my_films VALUES (
'{ "favorites" : [
  { "kind" : "comedy", "films" : [
    { "title" : "Bananas",
      "director" : "Woody Allen"},
    { "title" : "The Dinner Game",
      "director" : "Francis Veber" } ] },
  { "kind" : "horror", "films" : [
    { "title" : "Psycho",
      "director" : "Alfred Hitchcock" } ] },
  { "kind" : "thriller", "films" : [
    { "title" : "Vertigo",
      "director" : "Alfred Hitchcock" } ] },
  { "kind" : "drama", "films" : [
    { "title" : "Yojimbo",
      "director" : "Akira Kurosawa" } ] }
] }');
```

Query the `my_films` table holding some JSON data about the films and create a view that distributes the film genre, title, and director between separate columns:

```
SELECT jt.* FROM
my_films,
JSON_TABLE ( js, '$.favorites[*]' COLUMNS (
  id FOR ORDINALITY,
  kind text PATH '$.kind',
  NESTED PATH '$.films[*]' COLUMNS (
    title text PATH '$.title',
    director text PATH '$.director')) AS jt;
-----+-----+-----+-----
id |   kind   |      title      |   director
-----+-----+-----+-----
1  | comedy   | Bananas         | Woody Allen
1  | comedy   | The Dinner Game | Francis Veber
2  | horror    | Psycho          | Alfred Hitchcock
3  | thriller  | Vertigo         | Alfred Hitchcock
4  | drama     | Yojimbo         | Akira Kurosawa
(5 rows)
```

The following is a modified version of the above query to show the usage of `PASSING` arguments in the filter specified in the top-level JSON path expression and the various options for the individual columns:

```
SELECT jt.* FROM
my_films,
JSON_TABLE (js, '$.favorites[*] ? (@.films[*].director == $filter)'
```



```
PASSING 'Alfred Hitchcock' AS filter
  COLUMNS (
    id FOR ORDINALITY,
    kind text PATH '$.kind',
    title text FORMAT JSON PATH '$.films[*].title' OMIT QUOTES,
    director text PATH '$.films[*].director' KEEP QUOTES)) AS jt;
```

id	kind	title	director
1	horror	Psycho	"Alfred Hitchcock"
2	thriller	Vertigo	"Alfred Hitchcock"

(2 rows)

The following is a modified version of the above query to show the usage of `NESTED PATH` for populating title and director columns, illustrating how they are joined to the parent columns id and kind:

```
SELECT jt.* FROM
  my_films,
  JSON_TABLE ( js, '$.favorites[*] ? (@.films[*].director == $filter)'
    PASSING 'Alfred Hitchcock' AS filter
    COLUMNS (
      id FOR ORDINALITY,
      kind text PATH '$.kind',
      NESTED PATH '$.films[*]' COLUMNS (
        title text FORMAT JSON PATH '$.title' OMIT QUOTES,
        director text PATH '$.director' KEEP QUOTES))) AS jt;
```

id	kind	title	director
1	horror	Psycho	"Alfred Hitchcock"
2	thriller	Vertigo	"Alfred Hitchcock"

(2 rows)

The following is the same query but without the filter in the root path:

```
SELECT jt.* FROM
  my_films,
  JSON_TABLE ( js, '$.favorites[*]'
    COLUMNS (
      id FOR ORDINALITY,
      kind text PATH '$.kind',
      NESTED PATH '$.films[*]' COLUMNS (
        title text FORMAT JSON PATH '$.title' OMIT QUOTES,
        director text PATH '$.director' KEEP QUOTES))) AS jt;
```

id	kind	title	director
1	comedy	Bananas	"Woody Allen"
1	comedy	The Dinner Game	"Francis Veber"
2	horror	Psycho	"Alfred Hitchcock"
3	thriller	Vertigo	"Alfred Hitchcock"
4	drama	Yojimbo	"Akira Kurosawa"

(5 rows)

The following shows another query using a different `JSON` object as input. It shows the `UNION "sibling join"` between `NESTED` paths `$.movies[*]` and `$.books[*]` and also the usage of `FOR ORDINALITY` column at `NESTED` levels (columns `movie_id`, `book_id`, and `author_id`):

```
SELECT * FROM JSON_TABLE (
  '{"favorites":
    [{"movies":
      [{"name": "One", "director": "John Doe"}],
```

```

    {"name": "Two", "director": "Don Joe"}],
  "books":
    [{"name": "Mystery", "authors": [{"name": "Brown Dan"}]},
     {"name": "Wonder", "authors": [{"name": "Jun Murakami"}, {"name": "Craig Doe"}]}]
}}'::json, '$.favorites[*]'
COLUMNS (
  user_id FOR ORDINALITY,
  NESTED '$.movies[*]'
    COLUMNS (
      movie_id FOR ORDINALITY,
      mname text PATH '$.name',
      director text),
  NESTED '$.books[*]'
    COLUMNS (
      book_id FOR ORDINALITY,
      bname text PATH '$.name',
      NESTED '$.authors[*]'
        COLUMNS (
          author_id FOR ORDINALITY,
          author_name text PATH '$.name'))));

```

user_id	movie_id	mname	director	book_id	bname	author_id	author_name
1	1	One	John Doe				
1	2	Two	Don Joe				
1				1	Mystery	1	Brown Dan
1				2	Wonder	1	Jun Murakami
1				2	Wonder	2	Craig Doe

(5 rows)

9.17. Sequence Manipulation Functions

This section describes functions for operating on *sequence objects*, also called sequence generators or just sequences. Sequence objects are special single-row tables created with [CREATE SEQUENCE](#). Sequence objects are commonly used to generate unique identifiers for rows of a table. The sequence functions, listed in [Table 9.53](#), provide simple, multiuser-safe methods for obtaining successive sequence values from sequence objects.

Table 9.53. Sequence Functions

Function	Description
<code>nextval (regclass) → bigint</code>	Advances the sequence object to its next value and returns that value. This is done atomically: even if multiple sessions execute <code>nextval</code> concurrently, each will safely receive a distinct sequence value. If the sequence object has been created with default parameters, successive <code>nextval</code> calls will return successive values beginning with 1. Other behaviors can be obtained by using appropriate parameters in the CREATE SEQUENCE command. This function requires <code>USAGE</code> or <code>UPDATE</code> privilege on the sequence.
<code>setval (regclass, bigint [, boolean]) → bigint</code>	Sets the sequence object's current value, and optionally its <code>is_called</code> flag. The two-parameter form sets the sequence's <code>last_value</code> field to the specified value and sets its <code>is_called</code> field to <code>true</code> , meaning that the next <code>nextval</code> will advance the sequence before returning a value. The value that will be reported by <code>currval</code> is also set to the specified value. In the three-parameter form, <code>is_called</code> can be set to either <code>true</code> or <code>false</code> . <code>true</code> has the same effect as the two-parameter form. If it is set to <code>false</code> , the next <code>nextval</code> will return exactly the specified value, and sequence advancement commences with the following <code>nextval</code> . Furthermore, the value reported by <code>currval</code> is not changed in this case. For example,

Function	Description
<pre>SELECT setval('myseq', 42);</pre>	Next <code>nextval</code> will return 43
<pre>SELECT setval('myseq', 42, true);</pre>	Same as above
<pre>SELECT setval('myseq', 42, false);</pre>	Next <code>nextval</code> will return 42
<p>The result returned by <code>setval</code> is just the value of its second argument. This function requires <code>UPDATE</code> privilege on the sequence.</p>	
<pre>currval (regclass) → bigint</pre>	<p>Returns the value most recently obtained by <code>nextval</code> for this sequence in the current session. (An error is reported if <code>nextval</code> has never been called for this sequence in this session.) Because this is returning a session-local value, it gives a predictable answer whether or not other sessions have executed <code>nextval</code> since the current session did. This function requires <code>USAGE</code> or <code>SELECT</code> privilege on the sequence.</p>
<pre>lastval () → bigint</pre>	<p>Returns the value most recently returned by <code>nextval</code> in the current session. This function is identical to <code>currval</code>, except that instead of taking the sequence name as an argument it refers to whichever sequence <code>nextval</code> was most recently applied to in the current session. It is an error to call <code>lastval</code> if <code>nextval</code> has not yet been called in the current session. This function requires <code>USAGE</code> or <code>SELECT</code> privilege on the last used sequence.</p>

Caution

To avoid blocking concurrent transactions that obtain numbers from the same sequence, the value obtained by `nextval` is not reclaimed for re-use if the calling transaction later aborts. This means that transaction aborts or database crashes can result in gaps in the sequence of assigned values. That can happen without a transaction abort, too. For example an `INSERT` with an `ON CONFLICT` clause will compute the to-be-inserted tuple, including doing any required `nextval` calls, before detecting any conflict that would cause it to follow the `ON CONFLICT` rule instead. Thus, Postgres Pro sequence objects *cannot be used to obtain “gapless” sequences*.

Likewise, sequence state changes made by `setval` are immediately visible to other transactions, and are not undone if the calling transaction rolls back.

If the database cluster crashes before committing a transaction containing a `nextval` or `setval` call, the sequence state change might not have made its way to persistent storage, so that it is uncertain whether the sequence will have its original or updated state after the cluster restarts. This is harmless for usage of the sequence within the database, since other effects of uncommitted transactions will not be visible either. However, if you wish to use a sequence value for persistent outside-the-database purposes, make sure that the `nextval` call has been committed before doing so.

The sequence to be operated on by a sequence function is specified by a `regclass` argument, which is simply the OID of the sequence in the `pg_class` system catalog. You do not have to look up the OID by hand, however, since the `regclass` data type's input converter will do the work for you. See [Section 8.19](#) for details.

9.18. Conditional Expressions

This section describes the SQL-compliant conditional expressions available in Postgres Pro.

Tip

If your needs go beyond the capabilities of these conditional expressions, you might want to consider writing a server-side function in a more expressive programming language.

Note

Although `COALESCE`, `GREATEST`, and `LEAST` are syntactically similar to functions, they are not ordinary functions, and thus cannot be used with explicit `VARIADIC` array arguments.

9.18.1. CASE

The SQL `CASE` expression is a generic conditional expression, similar to `if/else` statements in other programming languages:

```
CASE WHEN condition THEN result
      [WHEN ...]
      [ELSE result]
END
```

`CASE` clauses can be used wherever an expression is valid. Each *condition* is an expression that returns a boolean result. If the condition's result is true, the value of the `CASE` expression is the *result* that follows the condition, and the remainder of the `CASE` expression is not processed. If the condition's result is not true, any subsequent `WHEN` clauses are examined in the same manner. If no `WHEN condition` yields true, the value of the `CASE` expression is the *result* of the `ELSE` clause. If the `ELSE` clause is omitted and no condition is true, the result is null.

An example:

```
SELECT * FROM test;
```

```
a
---
1
2
3
```

```
SELECT a,
       CASE WHEN a=1 THEN 'one'
            WHEN a=2 THEN 'two'
            ELSE 'other'
       END
FROM test;
```

```
a | case
---+-----
1 | one
2 | two
3 | other
```

The data types of all the *result* expressions must be convertible to a single output type. See [Section 10.5](#) for more details.

There is a “simple” form of `CASE` expression that is a variant of the general form above:

```
CASE expression
     WHEN value THEN result
     [WHEN ...]
     [ELSE result]
END
```

The first *expression* is computed, then compared to each of the *value* expressions in the `WHEN` clauses until one is found that is equal to it. If no match is found, the *result* of the `ELSE` clause (or a null value) is returned. This is similar to the `switch` statement in C.

The example above can be written using the simple CASE syntax:

```
SELECT a,
       CASE a WHEN 1 THEN 'one'
             WHEN 2 THEN 'two'
             ELSE 'other'
       END
FROM test;
```

```
a | case
---+-----
1 | one
2 | two
3 | other
```

A CASE expression does not evaluate any subexpressions that are not needed to determine the result. For example, this is a possible way of avoiding a division-by-zero failure:

```
SELECT ... WHERE CASE WHEN x <> 0 THEN y/x > 1.5 ELSE false END;
```

Note

As described in [Section 4.2.14](#), there are various situations in which subexpressions of an expression are evaluated at different times, so that the principle that “CASE evaluates only necessary subexpressions” is not ironclad. For example a constant `1/0` subexpression will usually result in a division-by-zero failure at planning time, even if it's within a CASE arm that would never be entered at run time.

9.18.2. COALESCE

```
COALESCE(value [, ...])
```

The COALESCE function returns the first of its arguments that is not null. Null is returned only if all arguments are null. It is often used to substitute a default value for null values when data is retrieved for display, for example:

```
SELECT COALESCE(description, short_description, '(none)') ...
```

This returns `description` if it is not null, otherwise `short_description` if it is not null, otherwise `(none)`.

The arguments must all be convertible to a common data type, which will be the type of the result (see [Section 10.5](#) for details).

Like a CASE expression, COALESCE only evaluates the arguments that are needed to determine the result; that is, arguments to the right of the first non-null argument are not evaluated. This SQL-standard function provides capabilities similar to NVL and IFNULL, which are used in some other database systems.

9.18.3. NULLIF

```
NULLIF(value1, value2)
```

The NULLIF function returns a null value if `value1` equals `value2`; otherwise it returns `value1`. This can be used to perform the inverse operation of the COALESCE example given above:

```
SELECT NULLIF(value, '(none)') ...
```

In this example, if `value` is `(none)`, null is returned, otherwise the value of `value` is returned.

The two arguments must be of comparable types. To be specific, they are compared exactly as if you had written `value1 = value2`, so there must be a suitable `=` operator available.

The result has the same type as the first argument — but there is a subtlety. What is actually returned is the first argument of the implied `=` operator, and in some cases that will have been promoted to match

the second argument's type. For example, `NULLIF(1, 2.2)` yields `numeric`, because there is no integer = numeric operator, only `numeric = numeric`.

9.18.4. GREATEST and LEAST

`GREATEST(value [, ...])`

`LEAST(value [, ...])`

The `GREATEST` and `LEAST` functions select the largest or smallest value from a list of any number of expressions. The expressions must all be convertible to a common data type, which will be the type of the result (see [Section 10.5](#) for details).

`NULL` values in the argument list are ignored. The result will be `NULL` only if all the expressions evaluate to `NULL`. (This is a deviation from the SQL standard. According to the standard, the return value is `NULL` if any argument is `NULL`. Some other databases behave this way.)

9.19. Array Functions and Operators

[Table 9.54](#) shows the specialized operators available for array types. In addition to those, the usual comparison operators shown in [Table 9.1](#) are available for arrays. The comparison operators compare the array contents element-by-element, using the default B-tree comparison function for the element data type, and sort based on the first difference. In multidimensional arrays the elements are visited in row-major order (last subscript varies most rapidly). If the contents of two arrays are equal but the dimensionality is different, the first difference in the dimensionality information determines the sort order.

Table 9.54. Array Operators

Operator	Description	Example(s)
<code>anyarray @> anyarray</code>	<code>→ boolean</code> Does the first array contain the second, that is, does each element appearing in the second array equal some element of the first array? (Duplicates are not treated specially, thus <code>ARRAY[1]</code> and <code>ARRAY[1,1]</code> are each considered to contain the other.)	<code>ARRAY[1,4,3] @> ARRAY[3,1,3] → t</code>
<code>anyarray <@ anyarray</code>	<code>→ boolean</code> Is the first array contained by the second?	<code>ARRAY[2,2,7] <@ ARRAY[1,7,4,2,6] → t</code>
<code>anyarray && anyarray</code>	<code>→ boolean</code> Do the arrays overlap, that is, have any elements in common?	<code>ARRAY[1,4,3] && ARRAY[2,1] → t</code>
<code>anycompatiblearray anycompatiblearray</code>	<code>→ anycompatiblearray</code> Concatenates the two arrays. Concatenating a null or empty array is a no-op; otherwise the arrays must have the same number of dimensions (as illustrated by the first example) or differ in number of dimensions by one (as illustrated by the second). If the arrays are not of identical element types, they will be coerced to a common type (see Section 10.5).	<code>ARRAY[1,2,3] ARRAY[4,5,6,7] → {1,2,3,4,5,6,7}</code> <code>ARRAY[1,2,3] ARRAY[[4,5,6],[7,8,9.9]] → {{1,2,3},{4,5,6},{7,8,9.9}}</code>
<code>anycompatible anycompatiblearray</code>	<code>→ anycompatiblearray</code> Concatenates an element onto the front of an array (which must be empty or one-dimensional).	<code>3 ARRAY[4,5,6] → {3,4,5,6}</code>
<code>anycompatiblearray anycompatible</code>	<code>→ anycompatiblearray</code>	

Operator	Description
Example(s)	Concatenates an element onto the end of an array (which must be empty or one-dimensional). <code>ARRAY[4,5,6] 7 → {4,5,6,7}</code>

See [Section 8.15](#) for more details about array operator behavior. See [Section 11.2](#) for more details about which operators support indexed operations.

[Table 9.55](#) shows the functions available for use with array types. See [Section 8.15](#) for more information and examples of the use of these functions.

Table 9.55. Array Functions

Function	Description
Example(s)	
<code>array_append (anycompatiblearray, anycompatible) → anycompatiblearray</code> Appends an element to the end of an array (same as the <code>anycompatiblearray anycompatible</code> operator). <code>array_append(ARRAY[1,2], 3) → {1,2,3}</code>	
<code>array_cat (anycompatiblearray, anycompatiblearray) → anycompatiblearray</code> Concatenates two arrays (same as the <code>anycompatiblearray anycompatiblearray</code> operator). <code>array_cat(ARRAY[1,2,3], ARRAY[4,5]) → {1,2,3,4,5}</code>	
<code>array_dims (anyarray) → text</code> Returns a text representation of the array's dimensions. <code>array_dims(ARRAY[[1,2,3], [4,5,6]]) → [1:2][1:3]</code>	
<code>array_fill (anyelement, integer[] [, integer[]]) → anyarray</code> Returns an array filled with copies of the given value, having dimensions of the lengths specified by the second argument. The optional third argument supplies lower-bound values for each dimension (which default to all 1). <code>array_fill(11, ARRAY[2,3]) → {{11,11,11},{11,11,11}}</code> <code>array_fill(7, ARRAY[3], ARRAY[2]) → [2:4]={7,7,7}</code>	
<code>array_length (anyarray, integer) → integer</code> Returns the length of the requested array dimension. (Produces NULL instead of 0 for empty or missing array dimensions.) <code>array_length(array[1,2,3], 1) → 3</code> <code>array_length(array[]::int[], 1) → NULL</code> <code>array_length(array['text'], 2) → NULL</code>	
<code>array_lower (anyarray, integer) → integer</code> Returns the lower bound of the requested array dimension. <code>array_lower('[0:2]={1,2,3}'::integer[], 1) → 0</code>	
<code>array_ndims (anyarray) → integer</code> Returns the number of dimensions of the array. <code>array_ndims(ARRAY[[1,2,3], [4,5,6]]) → 2</code>	
<code>array_position (anycompatiblearray, anycompatible [, integer]) → integer</code> Returns the subscript of the first occurrence of the second argument in the array, or NULL if it's not present. If the third argument is given, the search begins at that subscript. The array must be one-dimensional. Comparisons are done using IS NOT DISTINCT FROM semantics, so it is possible to search for NULL.	

Function	Description	Example(s)
		<code>array_position</code> (<code>ARRAY['sun', 'mon', 'tue', 'wed', 'thu', 'fri', 'sat'], 'mon'</code>) \rightarrow 2
	<code>array_positions</code> (<code>anycompatiblearray</code> , <code>anycompatible</code>) \rightarrow <code>integer[]</code> Returns an array of the subscripts of all occurrences of the second argument in the array given as first argument. The array must be one-dimensional. Comparisons are done using <code>IS NOT DISTINCT FROM</code> semantics, so it is possible to search for <code>NULL</code> . <code>NULL</code> is returned only if the array is <code>NULL</code> ; if the value is not found in the array, an empty array is returned.	<code>array_positions</code> (<code>ARRAY['A', 'A', 'B', 'A'], 'A'</code>) \rightarrow {1,2,4}
	<code>array_prepend</code> (<code>anycompatible</code> , <code>anycompatiblearray</code>) \rightarrow <code>anycompatiblearray</code> Prepends an element to the beginning of an array (same as the <code>anycompatible anycompatiblearray</code> operator).	<code>array_prepend</code> (1, <code>ARRAY[2,3]</code>) \rightarrow {1,2,3}
	<code>array_remove</code> (<code>anycompatiblearray</code> , <code>anycompatible</code>) \rightarrow <code>anycompatiblearray</code> Removes all elements equal to the given value from the array. The array must be one-dimensional. Comparisons are done using <code>IS NOT DISTINCT FROM</code> semantics, so it is possible to remove <code>NULL</code> s.	<code>array_remove</code> (<code>ARRAY[1,2,3,2]</code> , 2) \rightarrow {1,3}
	<code>array_replace</code> (<code>anycompatiblearray</code> , <code>anycompatible</code> , <code>anycompatible</code>) \rightarrow <code>anycompatiblearray</code> Replaces each array element equal to the second argument with the third argument.	<code>array_replace</code> (<code>ARRAY[1,2,5,4]</code> , 5, 3) \rightarrow {1,2,3,4}
	<code>array_sample</code> (<code>array anyarray</code> , <code>n integer</code>) \rightarrow <code>anyarray</code> Returns an array of <i>n</i> items randomly selected from <i>array</i> . <i>n</i> may not exceed the length of <i>array</i> 's first dimension. If <i>array</i> is multi-dimensional, an "item" is a slice having a given first subscript.	<code>array_sample</code> (<code>ARRAY[1,2,3,4,5,6]</code> , 3) \rightarrow {2,6,1} <code>array_sample</code> (<code>ARRAY[[1,2],[3,4],[5,6]]</code> , 2) \rightarrow {{5,6},{1,2}}
	<code>array_shuffle</code> (<code>anyarray</code>) \rightarrow <code>anyarray</code> Randomly shuffles the first dimension of the array.	<code>array_shuffle</code> (<code>ARRAY[[1,2],[3,4],[5,6]]</code>) \rightarrow {{5,6},{1,2},{3,4}}
	<code>array_to_string</code> (<code>array anyarray</code> , <code>delimiter text [, null_string text]</code>) \rightarrow <code>text</code> Converts each array element to its text representation, and concatenates those separated by the <i>delimiter</i> string. If <i>null_string</i> is given and is not <code>NULL</code> , then <code>NULL</code> array entries are represented by that string; otherwise, they are omitted. See also string_to_array .	<code>array_to_string</code> (<code>ARRAY[1, 2, 3, NULL, 5]</code> , ' ', '*') \rightarrow 1,2,3,*,5
	<code>array_upper</code> (<code>anyarray</code> , <code>integer</code>) \rightarrow <code>integer</code> Returns the upper bound of the requested array dimension.	<code>array_upper</code> (<code>ARRAY[1,8,3,7]</code> , 1) \rightarrow 4
	<code>cardinality</code> (<code>anyarray</code>) \rightarrow <code>integer</code> Returns the total number of elements in the array, or 0 if the array is empty.	<code>cardinality</code> (<code>ARRAY[[1,2],[3,4]]</code>) \rightarrow 4
	<code>trim_array</code> (<code>array anyarray</code> , <code>n integer</code>) \rightarrow <code>anyarray</code> Trims an array by removing the last <i>n</i> elements. If the array is multidimensional, only the first dimension is trimmed.	<code>trim_array</code> (<code>ARRAY[1,2,3,4,5,6]</code> , 2) \rightarrow {1,2,3,4}

Function	Description	Example(s)
<code>unnest (anyarray) → setof anyelement</code>	Expands an array into a set of rows. The array's elements are read out in storage order.	<pre>unnest (ARRAY[1,2]) → 1 2 unnest (ARRAY[['foo','bar'],['baz','quux']]) → foo bar baz quux</pre>
<code>unnest (anyarray, anyarray [, ...]) → setof anyelement, anyelement [, ...]</code>	Expands multiple arrays (possibly of different data types) into a set of rows. If the arrays are not all the same length then the shorter ones are padded with <code>NULLs</code> . This form is only allowed in a query's <code>FROM</code> clause; see Section 7.2.1.4 .	<pre>select * from unnest (ARRAY[1,2], ARRAY['foo','bar','baz']) as x(a,b) → a b ---+----- 1 foo 2 bar baz</pre>

See also [Section 9.21](#) about the aggregate function `array_agg` for use with arrays.

9.20. Range/Multirange Functions and Operators

See [Section 8.17](#) for an overview of range types.

[Table 9.56](#) shows the specialized operators available for range types. [Table 9.57](#) shows the specialized operators available for multirange types. In addition to those, the usual comparison operators shown in [Table 9.1](#) are available for range and multirange types. The comparison operators order first by the range lower bounds, and only if those are equal do they compare the upper bounds. The multirange operators compare each range until one is unequal. This does not usually result in a useful overall ordering, but the operators are provided to allow unique indexes to be constructed on ranges.

Table 9.56. Range Operators

Operator	Description	Example(s)
<code>anyrange @> anyrange → boolean</code>	Does the first range contain the second?	<pre>int4range(2,4) @> int4range(2,3) → t</pre>
<code>anyrange @> anyelement → boolean</code>	Does the range contain the element?	<pre>'[2011-01-01,2011-03-01) '::tsrange @> '2011-01-10' '::timestamp → t</pre>
<code>anyrange <@ anyrange → boolean</code>	Is the first range contained by the second?	<pre>int4range(2,4) <@ int4range(1,7) → t</pre>
<code>anyelement <@ anyrange → boolean</code>		

Operator	Description	Example(s)
	Is the element contained in the range?	<code>42 < @ int4range(1, 7) → f</code>
<code>anyrange && anyrange → boolean</code>	Do the ranges overlap, that is, have any elements in common?	<code>int8range(3, 7) && int8range(4, 12) → t</code>
<code>anyrange << anyrange → boolean</code>	Is the first range strictly left of the second?	<code>int8range(1, 10) << int8range(100, 110) → t</code>
<code>anyrange >> anyrange → boolean</code>	Is the first range strictly right of the second?	<code>int8range(50, 60) >> int8range(20, 30) → t</code>
<code>anyrange &< anyrange → boolean</code>	Does the first range not extend to the right of the second?	<code>int8range(1, 20) &< int8range(18, 20) → t</code>
<code>anyrange &> anyrange → boolean</code>	Does the first range not extend to the left of the second?	<code>int8range(7, 20) &> int8range(5, 10) → t</code>
<code>anyrange - - anyrange → boolean</code>	Are the ranges adjacent?	<code>numrange(1.1, 2.2) - - numrange(2.2, 3.3) → t</code>
<code>anyrange + anyrange → anyrange</code>	Computes the union of the ranges. The ranges must overlap or be adjacent, so that the union is a single range (but see <code>range_merge()</code>).	<code>numrange(5, 15) + numrange(10, 20) → [5, 20)</code>
<code>anyrange * anyrange → anyrange</code>	Computes the intersection of the ranges.	<code>int8range(5, 15) * int8range(10, 20) → [10, 15)</code>
<code>anyrange - anyrange → anyrange</code>	Computes the difference of the ranges. The second range must not be contained in the first in such a way that the difference would not be a single range.	<code>int8range(5, 15) - int8range(10, 20) → [5, 10)</code>

Table 9.57. Multirange Operators

Operator	Description	Example(s)
<code>anymultirange @> anymultirange → boolean</code>	Does the first multirange contain the second?	<code>'{[2, 4)]}'::int4multirange @> '[2, 3)]'::int4multirange → t</code>
<code>anymultirange @> anyrange → boolean</code>	Does the multirange contain the range?	<code>'{[2, 4)]}'::int4multirange @> int4range(2, 3) → t</code>
<code>anymultirange @> anyelement → boolean</code>	Does the multirange contain the element?	

Operator	Description	Example(s)
		<code>'{[2011-01-01,2011-03-01)}'::tsmultirange @> '2011-01-10'::timestamp → t</code>
<code>anyrange @> anymultirange → boolean</code>	Does the range contain the multirange?	<code>'[2,4)'::int4range @> '{[2,3)}'::int4multirange → t</code>
<code>anymultirange <@ anymultirange → boolean</code>	Is the first multirange contained by the second?	<code>'{[2,4)}'::int4multirange <@ '{[1,7)}'::int4multirange → t</code>
<code>anymultirange <@ anyrange → boolean</code>	Is the multirange contained by the range?	<code>'{[2,4)}'::int4multirange <@ int4range(1,7) → t</code>
<code>anyrange <@ anymultirange → boolean</code>	Is the range contained by the multirange?	<code>int4range(2,4) <@ '{[1,7)}'::int4multirange → t</code>
<code>anyelement <@ anymultirange → boolean</code>	Is the element contained by the multirange?	<code>4 <@ '{[1,7)}'::int4multirange → t</code>
<code>anymultirange && anymultirange → boolean</code>	Do the multiranges overlap, that is, have any elements in common?	<code>'{[3,7)}'::int8multirange && '{[4,12)}'::int8multirange → t</code>
<code>anymultirange && anyrange → boolean</code>	Does the multirange overlap the range?	<code>'{[3,7)}'::int8multirange && int8range(4,12) → t</code>
<code>anyrange && anymultirange → boolean</code>	Does the range overlap the multirange?	<code>int8range(3,7) && '{[4,12)}'::int8multirange → t</code>
<code>anymultirange << anymultirange → boolean</code>	Is the first multirange strictly left of the second?	<code>'{[1,10)}'::int8multirange << '{[100,110)}'::int8multirange → t</code>
<code>anymultirange << anyrange → boolean</code>	Is the multirange strictly left of the range?	<code>'{[1,10)}'::int8multirange << int8range(100,110) → t</code>
<code>anyrange << anymultirange → boolean</code>	Is the range strictly left of the multirange?	<code>int8range(1,10) << '{[100,110)}'::int8multirange → t</code>
<code>anymultirange >> anymultirange → boolean</code>	Is the first multirange strictly right of the second?	<code>'{[50,60)}'::int8multirange >> '{[20,30)}'::int8multirange → t</code>
<code>anymultirange >> anyrange → boolean</code>	Is the multirange strictly right of the range?	<code>'{[50,60)}'::int8multirange >> int8range(20,30) → t</code>
<code>anyrange >> anymultirange → boolean</code>	Is the range strictly right of the multirange?	

Operator	Description	Example(s)
		<code>int8range(50,60) >> '{[20,30)}'::int8multirange → t</code>
	<code>anymultirange &< anymultirange → boolean</code> Does the first multirange not extend to the right of the second?	<code>'{[1,20)}'::int8multirange &< '{[18,20)}'::int8multirange → t</code>
	<code>anymultirange &< anyrange → boolean</code> Does the multirange not extend to the right of the range?	<code>'{[1,20)}'::int8multirange &< int8range(18,20) → t</code>
	<code>anyrange &< anymultirange → boolean</code> Does the range not extend to the right of the multirange?	<code>int8range(1,20) &< '{[18,20)}'::int8multirange → t</code>
	<code>anymultirange &> anymultirange → boolean</code> Does the first multirange not extend to the left of the second?	<code>'{[7,20)}'::int8multirange &> '{[5,10)}'::int8multirange → t</code>
	<code>anymultirange &> anyrange → boolean</code> Does the multirange not extend to the left of the range?	<code>'{[7,20)}'::int8multirange &> int8range(5,10) → t</code>
	<code>anyrange &> anymultirange → boolean</code> Does the range not extend to the left of the multirange?	<code>int8range(7,20) &> '{[5,10)}'::int8multirange → t</code>
	<code>anymultirange - - anymultirange → boolean</code> Are the multiranges adjacent?	<code>'{[1.1,2.2)}'::nummultirange - - '{[2.2,3.3)}'::nummultirange → t</code>
	<code>anymultirange - - anyrange → boolean</code> Is the multirange adjacent to the range?	<code>'{[1.1,2.2)}'::nummultirange - - numrange(2.2,3.3) → t</code>
	<code>anyrange - - anymultirange → boolean</code> Is the range adjacent to the multirange?	<code>numrange(1.1,2.2) - - '{[2.2,3.3)}'::nummultirange → t</code>
	<code>anymultirange + anymultirange → anymultirange</code> Computes the union of the multiranges. The multiranges need not overlap or be adjacent.	<code>'{[5,10)}'::nummultirange + '{[15,20)}'::nummultirange → {[5,10), [15,20)}</code>
	<code>anymultirange * anymultirange → anymultirange</code> Computes the intersection of the multiranges.	<code>'{[5,15)}'::int8multirange * '{[10,20)}'::int8multirange → {[10,15)}</code>
	<code>anymultirange - anymultirange → anymultirange</code> Computes the difference of the multiranges.	<code>'{[5,20)}'::int8multirange - '{[10,15)}'::int8multirange → {[5,10), [15,20)}</code>

The left-of/right-of/adjacent operators always return false when an empty range or multirange is involved; that is, an empty range is not considered to be either before or after any other range.

Elsewhere empty ranges and multiranges are treated as the additive identity: anything unioned with an empty value is itself. Anything minus an empty value is itself. An empty multirange has exactly the

same points as an empty range. Every range contains the empty range. Every multirange contains as many empty ranges as you like.

The range union and difference operators will fail if the resulting range would need to contain two disjoint sub-ranges, as such a range cannot be represented. There are separate operators for union and difference that take multirange parameters and return a multirange, and they do not fail even if their arguments are disjoint. So if you need a union or difference operation for ranges that may be disjoint, you can avoid errors by first casting your ranges to multiranges.

[Table 9.58](#) shows the functions available for use with range types. [Table 9.59](#) shows the functions available for use with multirange types.

Table 9.58. Range Functions

Function	Description	Example(s)
<code>lower (anyrange)</code>	<code>→ anyelement</code> Extracts the lower bound of the range (NULL if the range is empty or has no lower bound).	<code>lower (numrange (1.1, 2.2))</code> <code>→ 1.1</code>
<code>upper (anyrange)</code>	<code>→ anyelement</code> Extracts the upper bound of the range (NULL if the range is empty or has no upper bound).	<code>upper (numrange (1.1, 2.2))</code> <code>→ 2.2</code>
<code>isempty (anyrange)</code>	<code>→ boolean</code> Is the range empty?	<code>isempty (numrange (1.1, 2.2))</code> <code>→ f</code>
<code>lower_inc (anyrange)</code>	<code>→ boolean</code> Is the range's lower bound inclusive?	<code>lower_inc (numrange (1.1, 2.2))</code> <code>→ t</code>
<code>upper_inc (anyrange)</code>	<code>→ boolean</code> Is the range's upper bound inclusive?	<code>upper_inc (numrange (1.1, 2.2))</code> <code>→ f</code>
<code>lower_inf (anyrange)</code>	<code>→ boolean</code> Does the range have no lower bound? (A lower bound of <code>-Infinity</code> returns false.)	<code>lower_inf ('(,)'::daterange)</code> <code>→ t</code>
<code>upper_inf (anyrange)</code>	<code>→ boolean</code> Does the range have no upper bound? (An upper bound of <code>Infinity</code> returns false.)	<code>upper_inf ('(,)'::daterange)</code> <code>→ t</code>
<code>range_merge (anyrange, anyrange)</code>	<code>→ anyrange</code> Computes the smallest range that includes both of the given ranges.	<code>range_merge ('[1,2)'::int4range, '[3,4)'::int4range)</code> <code>→ [1,4)</code>

Table 9.59. Multirange Functions

Function	Description	Example(s)
<code>lower (anymultirange)</code>	<code>→ anyelement</code> Extracts the lower bound of the multirange (NULL if the multirange is empty has no lower bound).	

Function	Description	Example(s)
		<code>lower('{[1.1,2.2)}'::nummultirange) → 1.1</code>
	<code>upper (anymultirange) → anyelement</code> Extracts the upper bound of the multirange (NULL if the multirange is empty or has no upper bound).	<code>upper('{[1.1,2.2)}'::nummultirange) → 2.2</code>
	<code>isempty (anymultirange) → boolean</code> Is the multirange empty?	<code>isempty('{[1.1,2.2)}'::nummultirange) → f</code>
	<code>lower_inc (anymultirange) → boolean</code> Is the multirange's lower bound inclusive?	<code>lower_inc('{[1.1,2.2)}'::nummultirange) → t</code>
	<code>upper_inc (anymultirange) → boolean</code> Is the multirange's upper bound inclusive?	<code>upper_inc('{[1.1,2.2)}'::nummultirange) → f</code>
	<code>lower_inf (anymultirange) → boolean</code> Does the multirange have no lower bound? (A lower bound of -Infinity returns false.)	<code>lower_inf('{(,)}'::datemultirange) → t</code>
	<code>upper_inf (anymultirange) → boolean</code> Does the multirange have no upper bound? (An upper bound of Infinity returns false.)	<code>upper_inf('{(,)}'::datemultirange) → t</code>
	<code>range_merge (anymultirange) → anyrange</code> Computes the smallest range that includes the entire multirange.	<code>range_merge('{[1,2), [3,4)}'::int4multirange) → [1,4)</code>
	<code>multirange (anyrange) → anymultirange</code> Returns a multirange containing just the given range.	<code>multirange('[1,2)'::int4range) → {[1,2)}</code>
	<code>unnest (anymultirange) → setof anyrange</code> Expands a multirange into a set of ranges in ascending order.	<code>unnest('{[1,2), [3,4)}'::int4multirange) →</code> <code>[1,2)</code> <code>[3,4)</code>

The `lower_inc`, `upper_inc`, `lower_inf`, and `upper_inf` functions all return false for an empty range or multirange.

9.21. Aggregate Functions

Aggregate functions compute a single result from a set of input values. The built-in general-purpose aggregate functions are listed in [Table 9.60](#) while statistical aggregates are in [Table 9.61](#). The built-in within-group ordered-set aggregate functions are listed in [Table 9.62](#) while the built-in within-group hypothetical-set ones are in [Table 9.63](#). Grouping operations, which are closely related to aggregate functions, are listed in [Table 9.64](#). The special syntax considerations for aggregate functions are explained in [Section 4.2.7](#). Consult [Section 2.7](#) for additional introductory information.

Aggregate functions that support *Partial Mode* are eligible to participate in various optimizations, such as parallel aggregation.

While all aggregates below accept an optional `ORDER BY` clause (as outlined in [Section 4.2.7](#)), the clause has only been added to aggregates whose output is affected by ordering.

Table 9.60. General-Purpose Aggregate Functions

Function Description	Partial Mode
<code>any_value (anyelement) → same as input type</code> Returns an arbitrary value from the non-null input values.	Yes
<code>array_agg (anynonarray ORDER BY input_sort_columns) → anyarray</code> Collects all the input values, including nulls, into an array.	Yes
<code>array_agg (anyarray ORDER BY input_sort_columns) → anyarray</code> Concatenates all the input arrays into an array of one higher dimension. (The inputs must all have the same dimensionality, and cannot be empty or null.)	Yes
<code>avg (smallint) → numeric</code> <code>avg (integer) → numeric</code> <code>avg (bigint) → numeric</code> <code>avg (numeric) → numeric</code> <code>avg (real) → double precision</code> <code>avg (double precision) → double precision</code> <code>avg (interval) → interval</code> Computes the average (arithmetic mean) of all the non-null input values.	Yes
<code>bit_and (smallint) → smallint</code> <code>bit_and (integer) → integer</code> <code>bit_and (bigint) → bigint</code> <code>bit_and (bit) → bit</code> Computes the bitwise AND of all non-null input values.	Yes
<code>bit_or (smallint) → smallint</code> <code>bit_or (integer) → integer</code> <code>bit_or (bigint) → bigint</code> <code>bit_or (bit) → bit</code> Computes the bitwise OR of all non-null input values.	Yes
<code>bit_xor (smallint) → smallint</code> <code>bit_xor (integer) → integer</code> <code>bit_xor (bigint) → bigint</code> <code>bit_xor (bit) → bit</code> Computes the bitwise exclusive OR of all non-null input values. Can be useful as a checksum for an unordered set of values.	Yes
<code>bool_and (boolean) → boolean</code> Returns true if all non-null input values are true, otherwise false.	Yes
<code>bool_or (boolean) → boolean</code> Returns true if any non-null input value is true, otherwise false.	Yes
<code>count (*) → bigint</code> Computes the number of input rows.	Yes
<code>count ("any") → bigint</code> Computes the number of input rows in which the input value is not null.	Yes
<code>every (boolean) → boolean</code>	Yes

Function Description	Partial Mode
This is the SQL standard's equivalent to <code>bool_and</code> .	
<code>json_agg (anyelement ORDER BY input_sort_columns) → json</code> <code>jsonb_agg (anyelement ORDER BY input_sort_columns) → jsonb</code> Collects all the input values, including nulls, into a JSON array. Values are converted to JSON as per <code>to_json</code> or <code>to_jsonb</code> .	No
<code>json_agg_strict (anyelement) → json</code> <code>jsonb_agg_strict (anyelement) → jsonb</code> Collects all the input values, skipping nulls, into a JSON array. Values are converted to JSON as per <code>to_json</code> or <code>to_jsonb</code> .	No
<code>json_arrayagg ([value_expression] [ORDER BY sort_expression] [{ NULL ABSENT } ON NULL] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</code> Behaves in the same way as <code>json_array</code> but as an aggregate function so it only takes one <code>value_expression</code> parameter. If <code>ABSENT ON NULL</code> is specified, any NULL values are omitted. If <code>ORDER BY</code> is specified, the elements will appear in the array in that order rather than in the input order. <code>SELECT json_arrayagg(v) FROM (VALUES(2),(1)) t(v) → [2, 1]</code>	No
<code>json_objectagg ([{ key_expression { VALUE ':' } value_expression }] [{ NULL ABSENT } ON NULL] [{ WITH WITHOUT } UNIQUE [KEYS]] [RETURNING data_type [FORMAT JSON [ENCODING UTF8]]])</code> Behaves like <code>json_object</code> , but as an aggregate function, so it only takes one <code>key_expression</code> and one <code>value_expression</code> parameter. <code>SELECT json_objectagg(k:v) FROM (VALUES ('a'::text,current_date), ('b',current_date + 1)) AS t(k,v) → { "a" : "2022-05-10", "b" : "2022-05-11" }</code>	No
<code>json_object_agg (key "any", value "any" ORDER BY input_sort_columns) → json</code> <code>jsonb_object_agg (key "any", value "any" ORDER BY input_sort_columns) → jsonb</code> Collects all the key/value pairs into a JSON object. Key arguments are coerced to text; value arguments are converted as per <code>to_json</code> or <code>to_jsonb</code> . Values can be null, but keys cannot.	No
<code>json_object_agg_strict (key "any", value "any") → json</code> <code>jsonb_object_agg_strict (key "any", value "any") → jsonb</code> Collects all the key/value pairs into a JSON object. Key arguments are coerced to text; value arguments are converted as per <code>to_json</code> or <code>to_jsonb</code> . The <code>key</code> can not be null. If the <code>value</code> is null then the entry is skipped,	No
<code>json_object_agg_unique (key "any", value "any") → json</code> <code>jsonb_object_agg_unique (key "any", value "any") → jsonb</code> Collects all the key/value pairs into a JSON object. Key arguments are coerced to text; value arguments are converted as per <code>to_json</code> or <code>to_jsonb</code> . Values can be null, but keys cannot. If there is a duplicate key an error is thrown.	No
<code>json_object_agg_unique_strict (key "any", value "any") → json</code> <code>jsonb_object_agg_unique_strict (key "any", value "any") → jsonb</code> Collects all the key/value pairs into a JSON object. Key arguments are coerced to text; value arguments are converted as per <code>to_json</code> or <code>to_jsonb</code> . The <code>key</code> can not be null. If the <code>value</code> is null then the entry is skipped. If there is a duplicate key an error is thrown.	No
<code>max (see text) → same as input type</code>	Yes

Function Description	Partial Mode
Computes the maximum of the non-null input values. Available for any numeric, string, date/time, or enum type, as well as inet, interval, money, oid, pg_lsn , tid, xid8, and arrays of any of these types.	
min (<i>see text</i>) → <i>same as input type</i> Computes the minimum of the non-null input values. Available for any numeric, string, date/time, or enum type, as well as inet, interval, money, oid, pg_lsn , tid, xid8, and arrays of any of these types.	Yes
range_agg (<i>value</i> anyrange) → anymultirange range_agg (<i>value</i> anymultirange) → anymultirange Computes the union of the non-null input values.	No
range_intersect_agg (<i>value</i> anyrange) → anyrange range_intersect_agg (<i>value</i> anymultirange) → anymultirange Computes the intersection of the non-null input values.	No
string_agg (<i>value</i> text, <i>delimiter</i> text) → text string_agg (<i>value</i> bytea, <i>delimiter</i> bytea ORDER BY input_sort_columns) → bytea Concatenates the non-null input values into a string. Each value after the first is preceded by the corresponding <i>delimiter</i> (if it's not null).	Yes
sum (smallint) → bigint sum (integer) → bigint sum (bigint) → numeric sum (numeric) → numeric sum (real) → real sum (double precision) → double precision sum (interval) → interval sum (money) → money Computes the sum of the non-null input values.	Yes
xmlagg (xml ORDER BY input_sort_columns) → xml Concatenates the non-null XML input values (see Section 9.15.1.8).	No

It should be noted that except for `count`, these functions return a null value when no rows are selected. In particular, `sum` of no rows returns null, not zero as one might expect, and `array_agg` returns null rather than an empty array when there are no input rows. The `coalesce` function can be used to substitute zero or an empty array for null when necessary.

The aggregate functions `array_agg`, `json_agg`, `jsonb_agg`, `json_agg_strict`, `jsonb_agg_strict`, `json_object_agg`, `jsonb_object_agg`, `json_object_agg_strict`, `jsonb_object_agg_strict`, `json_object_agg_unique`, `jsonb_object_agg_unique`, `json_object_agg_unique_strict`, `jsonb_object_agg_unique_strict`, `string_agg`, and `xmlagg`, as well as similar user-defined aggregate functions, produce meaningfully different result values depending on the order of the input values. This ordering is unspecified by default, but can be controlled by writing an `ORDER BY` clause within the aggregate call, as shown in [Section 4.2.7](#). Alternatively, supplying the input values from a sorted subquery will usually work. For example:

```
SELECT xmlagg(x) FROM (SELECT x FROM test ORDER BY y DESC) AS tab;
```

Beware that this approach can fail if the outer query level contains additional processing, such as a join, because that might cause the subquery's output to be reordered before the aggregate is computed.

Note

The boolean aggregates `bool_and` and `bool_or` correspond to the standard SQL aggregates `every` and `any` or `some`. Postgres Pro supports `every`, but not `any` or `some`, because there is an ambiguity built into the standard syntax:

```
SELECT b1 = ANY((SELECT b2 FROM t2 ...)) FROM t1 ...;
```

Here `ANY` can be considered either as introducing a subquery, or as being an aggregate function, if the subquery returns one row with a Boolean value. Thus the standard name cannot be given to these aggregates.

Note

Users accustomed to working with other SQL database management systems might be disappointed by the performance of the `count` aggregate when it is applied to the entire table. A query like:

```
SELECT count(*) FROM sometable;
```

will require effort proportional to the size of the table: Postgres Pro will need to scan either the entire table or the entirety of an index that includes all rows in the table.

Table 9.61 shows aggregate functions typically used in statistical analysis. (These are separated out merely to avoid cluttering the listing of more-commonly-used aggregates.) Functions shown as accepting *numeric_type* are available for all the types `smallint`, `integer`, `bigint`, `numeric`, `real`, and `double precision`. Where the description mentions *N*, it means the number of input rows for which all the input expressions are non-null. In all cases, null is returned if the computation is meaningless, for example when *N* is zero.

Table 9.61. Aggregate Functions for Statistics

Function Description	Partial Mode
<code>corr (Y double precision, X double precision) → double precision</code> Computes the correlation coefficient.	Yes
<code>covar_pop (Y double precision, X double precision) → double precision</code> Computes the population covariance.	Yes
<code>covar_samp (Y double precision, X double precision) → double precision</code> Computes the sample covariance.	Yes
<code>regr_avgx (Y double precision, X double precision) → double precision</code> Computes the average of the independent variable, $\text{sum}(X) / N$.	Yes
<code>regr_avgy (Y double precision, X double precision) → double precision</code> Computes the average of the dependent variable, $\text{sum}(Y) / N$.	Yes
<code>regr_count (Y double precision, X double precision) → bigint</code> Computes the number of rows in which both inputs are non-null.	Yes
<code>regr_intercept (Y double precision, X double precision) → double precision</code> Computes the y-intercept of the least-squares-fit linear equation determined by the (<i>X</i> , <i>y</i>) pairs.	Yes
<code>regr_r2 (Y double precision, X double precision) → double precision</code> Computes the square of the correlation coefficient.	Yes
<code>regr_slope (Y double precision, X double precision) → double precision</code>	Yes

Function Description	Partial Mode
Computes the slope of the least-squares-fit linear equation determined by the (X, Y) pairs.	
<code>regr_sxx (Y double precision, X double precision) → double precision</code> Computes the “sum of squares” of the independent variable, $\text{sum}(X^2) - \text{sum}(X)^2 / N$.	Yes
<code>regr_sxy (Y double precision, X double precision) → double precision</code> Computes the “sum of products” of independent times dependent variables, $\text{sum}(X * Y) - \text{sum}(X) * \text{sum}(Y) / N$.	Yes
<code>regr_syy (Y double precision, X double precision) → double precision</code> Computes the “sum of squares” of the dependent variable, $\text{sum}(Y^2) - \text{sum}(Y)^2 / N$.	Yes
<code>stddev (numeric_type) → double precision</code> for real or double precision, otherwise numeric This is a historical alias for <code>stddev_samp</code> .	Yes
<code>stddev_pop (numeric_type) → double precision</code> for real or double precision, otherwise numeric Computes the population standard deviation of the input values.	Yes
<code>stddev_samp (numeric_type) → double precision</code> for real or double precision, otherwise numeric Computes the sample standard deviation of the input values.	Yes
<code>variance (numeric_type) → double precision</code> for real or double precision, otherwise numeric This is a historical alias for <code>var_samp</code> .	Yes
<code>var_pop (numeric_type) → double precision</code> for real or double precision, otherwise numeric Computes the population variance of the input values (square of the population standard deviation).	Yes
<code>var_samp (numeric_type) → double precision</code> for real or double precision, otherwise numeric Computes the sample variance of the input values (square of the sample standard deviation).	Yes

Table 9.62 shows some aggregate functions that use the *ordered-set aggregate* syntax. These functions are sometimes referred to as “inverse distribution” functions. Their aggregated input is introduced by `ORDER BY`, and they may also take a *direct argument* that is not aggregated, but is computed only once. All these functions ignore null values in their aggregated input. For those that take a *fraction* parameter, the fraction value must be between 0 and 1; an error is thrown if not. However, a null *fraction* value simply produces a null result.

Table 9.62. Ordered-Set Aggregate Functions

Function Description	Partial Mode
<code>mode () WITHIN GROUP (ORDER BY anyelement) → anyelement</code> Computes the <i>mode</i> , the most frequent value of the aggregated argument (arbitrarily choosing the first one if there are multiple equally-frequent values). The aggregated argument must be of a sortable type.	No

Function Description	Partial Mode
<code>percentile_cont (fraction double precision) WITHIN GROUP (ORDER BY double precision) → double precision</code> <code>percentile_cont (fraction double precision) WITHIN GROUP (ORDER BY interval) → interval</code> Computes the <i>continuous percentile</i> , a value corresponding to the specified <i>fraction</i> within the ordered set of aggregated argument values. This will interpolate between adjacent input items if needed.	No
<code>percentile_cont (fractions double precision[]) WITHIN GROUP (ORDER BY double precision) → double precision[]</code> <code>percentile_cont (fractions double precision[]) WITHIN GROUP (ORDER BY interval) → interval[]</code> Computes multiple continuous percentiles. The result is an array of the same dimensions as the <i>fractions</i> parameter, with each non-null element replaced by the (possibly interpolated) value corresponding to that percentile.	No
<code>percentile_disc (fraction double precision) WITHIN GROUP (ORDER BY anyelement) → anyelement</code> Computes the <i>discrete percentile</i> , the first value within the ordered set of aggregated argument values whose position in the ordering equals or exceeds the specified <i>fraction</i> . The aggregated argument must be of a sortable type.	No
<code>percentile_disc (fractions double precision[]) WITHIN GROUP (ORDER BY anyelement) → anyarray</code> Computes multiple discrete percentiles. The result is an array of the same dimensions as the <i>fractions</i> parameter, with each non-null element replaced by the input value corresponding to that percentile. The aggregated argument must be of a sortable type.	No

Each of the “hypothetical-set” aggregates listed in [Table 9.63](#) is associated with a window function of the same name defined in [Section 9.22](#). In each case, the aggregate's result is the value that the associated window function would have returned for the “hypothetical” row constructed from *args*, if such a row had been added to the sorted group of rows represented by the *sorted_args*. For each of these functions, the list of direct arguments given in *args* must match the number and types of the aggregated arguments given in *sorted_args*. Unlike most built-in aggregates, these aggregates are not strict, that is they do not drop input rows containing nulls. Null values sort according to the rule specified in the `ORDER BY` clause.

Table 9.63. Hypothetical-Set Aggregate Functions

Function Description	Partial Mode
<code>rank (args) WITHIN GROUP (ORDER BY sorted_args) → bigint</code> Computes the rank of the hypothetical row, with gaps; that is, the row number of the first row in its peer group.	No
<code>dense_rank (args) WITHIN GROUP (ORDER BY sorted_args) → bigint</code> Computes the rank of the hypothetical row, without gaps; this function effectively counts peer groups.	No
<code>percent_rank (args) WITHIN GROUP (ORDER BY sorted_args) → double precision</code> Computes the relative rank of the hypothetical row, that is $(rank - 1) / (total\ rows - 1)$. The value thus ranges from 0 to 1 inclusive.	No
<code>cume_dist (args) WITHIN GROUP (ORDER BY sorted_args) → double precision</code> Computes the cumulative distribution, that is $(number\ of\ rows\ preceding\ or\ peers\ with\ hypothetical\ row) / (total\ rows)$. The value thus ranges from $1/N$ to 1.	No

Table 9.64. Grouping Operations

Function	Description
<code>GROUPING (<i>group_by_expression(s)</i>) → integer</code>	Returns a bit mask indicating which <code>GROUP BY</code> expressions are not included in the current grouping set. Bits are assigned with the rightmost argument corresponding to the least-significant bit; each bit is 0 if the corresponding expression is included in the grouping criteria of the grouping set generating the current result row, and 1 if it is not included.

The grouping operations shown in [Table 9.64](#) are used in conjunction with grouping sets (see [Section 7.2.4](#)) to distinguish result rows. The arguments to the `GROUPING` function are not actually evaluated, but they must exactly match expressions given in the `GROUP BY` clause of the associated query level. For example:

```
=> SELECT * FROM items_sold;
```

make	model	sales
-----+-----+-----		
Foo	GT	10
Foo	Tour	20
Bar	City	15
Bar	Sport	5

(4 rows)

```
=> SELECT make, model, GROUPING(make,model), sum(sales) FROM items_sold GROUP BY
ROLLUP (make,model);
```

make	model	grouping	sum
-----+-----+-----+-----			
Foo	GT	0	10
Foo	Tour	0	20
Bar	City	0	15
Bar	Sport	0	5
Foo		1	30
Bar		1	20
		3	50

(7 rows)

Here, the `grouping` value 0 in the first four rows shows that those have been grouped normally, over both the grouping columns. The value 1 indicates that `model` was not grouped by in the next-to-last two rows, and the value 3 indicates that neither `make` nor `model` was grouped by in the last row (which therefore is an aggregate over all the input rows).

9.22. Window Functions

Window functions provide the ability to perform calculations across sets of rows that are related to the current query row. See [Section 3.5](#) for an introduction to this feature, and [Section 4.2.8](#) for syntax details.

The built-in window functions are listed in [Table 9.65](#). Note that these functions *must* be invoked using window function syntax, i.e., an `OVER` clause is required.

In addition to these functions, any built-in or user-defined ordinary aggregate (i.e., not ordered-set or hypothetical-set aggregates) can be used as a window function; see [Section 9.21](#) for a list of the built-in aggregates. Aggregate functions act as window functions only when an `OVER` clause follows the call; otherwise they act as plain aggregates and return a single row for the entire set.

Table 9.65. General-Purpose Window Functions

Function	Description
<code>row_number () → bigint</code>	

Function	Description
	Returns the number of the current row within its partition, counting from 1.
<code>rank () → bigint</code>	Returns the rank of the current row, with gaps; that is, the <code>row_number</code> of the first row in its peer group.
<code>dense_rank () → bigint</code>	Returns the rank of the current row, without gaps; this function effectively counts peer groups.
<code>percent_rank () → double precision</code>	Returns the relative rank of the current row, that is $(\text{rank} - 1) / (\text{total partition rows} - 1)$. The value thus ranges from 0 to 1 inclusive.
<code>cume_dist () → double precision</code>	Returns the cumulative distribution, that is $(\text{number of partition rows preceding or peers with current row}) / (\text{total partition rows})$. The value thus ranges from $1/N$ to 1.
<code>ntile (num_buckets integer) → integer</code>	Returns an integer ranging from 1 to the argument value, dividing the partition as equally as possible.
<code>lag (value anycompatible [, offset integer [, default anycompatible]]) → anycompatible</code>	Returns <i>value</i> evaluated at the row that is <i>offset</i> rows before the current row within the partition; if there is no such row, instead returns <i>default</i> (which must be of a type compatible with <i>value</i>). Both <i>offset</i> and <i>default</i> are evaluated with respect to the current row. If omitted, <i>offset</i> defaults to 1 and <i>default</i> to NULL.
<code>lead (value anycompatible [, offset integer [, default anycompatible]]) → anycompatible</code>	Returns <i>value</i> evaluated at the row that is <i>offset</i> rows after the current row within the partition; if there is no such row, instead returns <i>default</i> (which must be of a type compatible with <i>value</i>). Both <i>offset</i> and <i>default</i> are evaluated with respect to the current row. If omitted, <i>offset</i> defaults to 1 and <i>default</i> to NULL.
<code>first_value (value anyelement) → anyelement</code>	Returns <i>value</i> evaluated at the row that is the first row of the window frame.
<code>last_value (value anyelement) → anyelement</code>	Returns <i>value</i> evaluated at the row that is the last row of the window frame.
<code>nth_value (value anyelement, n integer) → anyelement</code>	Returns <i>value</i> evaluated at the row that is the <i>n</i> 'th row of the window frame (counting from 1); returns NULL if there is no such row.

All of the functions listed in [Table 9.65](#) depend on the sort ordering specified by the `ORDER BY` clause of the associated window definition. Rows that are not distinct when considering only the `ORDER BY` columns are said to be *peers*. The four ranking functions (including `cume_dist`) are defined so that they give the same answer for all rows of a peer group.

Note that `first_value`, `last_value`, and `nth_value` consider only the rows within the “window frame”, which by default contains the rows from the start of the partition through the last peer of the current row. This is likely to give unhelpful results for `last_value` and sometimes also `nth_value`. You can redefine the frame by adding a suitable frame specification (`RANGE`, `ROWS` or `GROUPS`) to the `OVER` clause. See [Section 4.2.8](#) for more information about frame specifications.

When an aggregate function is used as a window function, it aggregates over the rows within the current row's window frame. An aggregate used with `ORDER BY` and the default window frame definition produces a “running sum” type of behavior, which may or may not be what's wanted. To obtain aggregation

over the whole partition, omit `ORDER BY` or use `ROWS BETWEEN UNBOUNDED PRECEDING AND UNBOUNDED FOLLOWING`. Other frame specifications can be used to obtain other effects.

Note

The SQL standard defines a `RESPECT NULLS` or `IGNORE NULLS` option for `lead`, `lag`, `first_value`, `last_value`, and `nth_value`. This is not implemented in Postgres Pro: the behavior is always the same as the standard's default, namely `RESPECT NULLS`. Likewise, the standard's `FROM FIRST` or `FROM LAST` option for `nth_value` is not implemented: only the default `FROM FIRST` behavior is supported. (You can achieve the result of `FROM LAST` by reversing the `ORDER BY` ordering.)

9.23. Merge Support Functions

Postgres Pro includes one merge support function that may be used in the `RETURNING` list of a [MERGE](#) command to identify the action taken for each row; see [Table 9.66](#).

Table 9.66. Merge Support Functions

Function	Description
<code>merge_action ()</code> → text	Returns the merge action command executed for the current row. This will be 'INSERT', 'UPDATE', or 'DELETE'.

Example:

```
MERGE INTO products p
  USING stock s ON p.product_id = s.product_id
  WHEN MATCHED AND s.quantity > 0 THEN
    UPDATE SET in_stock = true, quantity = s.quantity
  WHEN MATCHED THEN
    UPDATE SET in_stock = false, quantity = 0
  WHEN NOT MATCHED THEN
    INSERT (product_id, in_stock, quantity)
      VALUES (s.product_id, true, s.quantity)
  RETURNING merge_action(), p.*;
```

merge_action	product_id	in_stock	quantity
UPDATE	1001	t	50
UPDATE	1002	f	0
INSERT	1003	t	10

Note that this function can only be used in the `RETURNING` list of a `MERGE` command. It is an error to use it in any other part of a query.

9.24. Subquery Expressions

This section describes the SQL-compliant subquery expressions available in Postgres Pro. All of the expression forms documented in this section return Boolean (true/false) results.

9.24.1. EXISTS

`EXISTS (subquery)`

The argument of `EXISTS` is an arbitrary `SELECT` statement, or *subquery*. The subquery is evaluated to determine whether it returns any rows. If it returns at least one row, the result of `EXISTS` is “true”; if the subquery returns no rows, the result of `EXISTS` is “false”.

The subquery can refer to variables from the surrounding query, which will act as constants during any one evaluation of the subquery.

The subquery will generally only be executed long enough to determine whether at least one row is returned, not all the way to completion. It is unwise to write a subquery that has side effects (such as calling sequence functions); whether the side effects occur might be unpredictable.

Since the result depends only on whether any rows are returned, and not on the contents of those rows, the output list of the subquery is normally unimportant. A common coding convention is to write all `EXISTS` tests in the form `EXISTS (SELECT 1 WHERE ...)`. There are exceptions to this rule however, such as subqueries that use `INTERSECT`.

This simple example is like an inner join on `col2`, but it produces at most one output row for each `tab1` row, even if there are several matching `tab2` rows:

```
SELECT col1
FROM tab1
WHERE EXISTS (SELECT 1 FROM tab2 WHERE col2 = tab1.col2);
```

9.24.2. `IN`

expression `IN` (*subquery*)

The right-hand side is a parenthesized subquery, which must return exactly one column. The left-hand expression is evaluated and compared to each row of the subquery result. The result of `IN` is “true” if any equal subquery row is found. The result is “false” if no equal row is found (including the case where the subquery returns no rows).

Note that if the left-hand expression yields null, or if there are no equal right-hand values and at least one right-hand row yields null, the result of the `IN` construct will be null, not false. This is in accordance with SQL's normal rules for Boolean combinations of null values.

As with `EXISTS`, it's unwise to assume that the subquery will be evaluated completely.

row_constructor `IN` (*subquery*)

The left-hand side of this form of `IN` is a row constructor, as described in [Section 4.2.13](#). The right-hand side is a parenthesized subquery, which must return exactly as many columns as there are expressions in the left-hand row. The left-hand expressions are evaluated and compared row-wise to each row of the subquery result. The result of `IN` is “true” if any equal subquery row is found. The result is “false” if no equal row is found (including the case where the subquery returns no rows).

As usual, null values in the rows are combined per the normal rules of SQL Boolean expressions. Two rows are considered equal if all their corresponding members are non-null and equal; the rows are unequal if any corresponding members are non-null and unequal; otherwise the result of that row comparison is unknown (null). If all the per-row results are either unequal or null, with at least one null, then the result of `IN` is null.

9.24.3. `NOT IN`

expression `NOT IN` (*subquery*)

The right-hand side is a parenthesized subquery, which must return exactly one column. The left-hand expression is evaluated and compared to each row of the subquery result. The result of `NOT IN` is “true” if only unequal subquery rows are found (including the case where the subquery returns no rows). The result is “false” if any equal row is found.

Note that if the left-hand expression yields null, or if there are no equal right-hand values and at least one right-hand row yields null, the result of the `NOT IN` construct will be null, not true. This is in accordance with SQL's normal rules for Boolean combinations of null values.

As with `EXISTS`, it's unwise to assume that the subquery will be evaluated completely.

```
row_constructor NOT IN (subquery)
```

The left-hand side of this form of `NOT IN` is a row constructor, as described in [Section 4.2.13](#). The right-hand side is a parenthesized subquery, which must return exactly as many columns as there are expressions in the left-hand row. The left-hand expressions are evaluated and compared row-wise to each row of the subquery result. The result of `NOT IN` is “true” if only unequal subquery rows are found (including the case where the subquery returns no rows). The result is “false” if any equal row is found.

As usual, null values in the rows are combined per the normal rules of SQL Boolean expressions. Two rows are considered equal if all their corresponding members are non-null and equal; the rows are unequal if any corresponding members are non-null and unequal; otherwise the result of that row comparison is unknown (null). If all the per-row results are either unequal or null, with at least one null, then the result of `NOT IN` is null.

9.24.4. ANY/SOME

```
expression operator ANY (subquery)
expression operator SOME (subquery)
```

The right-hand side is a parenthesized subquery, which must return exactly one column. The left-hand expression is evaluated and compared to each row of the subquery result using the given *operator*, which must yield a Boolean result. The result of `ANY` is “true” if any true result is obtained. The result is “false” if no true result is found (including the case where the subquery returns no rows).

`SOME` is a synonym for `ANY`. `IN` is equivalent to `= ANY`.

Note that if there are no successes and at least one right-hand row yields null for the operator's result, the result of the `ANY` construct will be null, not false. This is in accordance with SQL's normal rules for Boolean combinations of null values.

As with `EXISTS`, it's unwise to assume that the subquery will be evaluated completely.

```
row_constructor operator ANY (subquery)
row_constructor operator SOME (subquery)
```

The left-hand side of this form of `ANY` is a row constructor, as described in [Section 4.2.13](#). The right-hand side is a parenthesized subquery, which must return exactly as many columns as there are expressions in the left-hand row. The left-hand expressions are evaluated and compared row-wise to each row of the subquery result, using the given *operator*. The result of `ANY` is “true” if the comparison returns true for any subquery row. The result is “false” if the comparison returns false for every subquery row (including the case where the subquery returns no rows). The result is `NULL` if no comparison with a subquery row returns true, and at least one comparison returns `NULL`.

See [Section 9.25.5](#) for details about the meaning of a row constructor comparison.

9.24.5. ALL

```
expression operator ALL (subquery)
```

The right-hand side is a parenthesized subquery, which must return exactly one column. The left-hand expression is evaluated and compared to each row of the subquery result using the given *operator*, which must yield a Boolean result. The result of `ALL` is “true” if all rows yield true (including the case where the subquery returns no rows). The result is “false” if any false result is found. The result is `NULL` if no comparison with a subquery row returns false, and at least one comparison returns `NULL`.

`NOT IN` is equivalent to `<> ALL`.

As with `EXISTS`, it's unwise to assume that the subquery will be evaluated completely.

```
row_constructor operator ALL (subquery)
```

The left-hand side of this form of `ALL` is a row constructor, as described in [Section 4.2.13](#). The right-hand side is a parenthesized subquery, which must return exactly as many columns as there are expressions in the left-hand row. The left-hand expressions are evaluated and compared row-wise to each row of the subquery result, using the given *operator*. The result of `ALL` is “true” if the comparison returns true for all subquery rows (including the case where the subquery returns no rows). The result is “false” if the comparison returns false for any subquery row. The result is `NULL` if no comparison with a subquery row returns false, and at least one comparison returns `NULL`.

See [Section 9.25.5](#) for details about the meaning of a row constructor comparison.

9.24.6. Single-Row Comparison

```
row_constructor operator (subquery)
```

The left-hand side is a row constructor, as described in [Section 4.2.13](#). The right-hand side is a parenthesized subquery, which must return exactly as many columns as there are expressions in the left-hand row. Furthermore, the subquery cannot return more than one row. (If it returns zero rows, the result is taken to be null.) The left-hand side is evaluated and compared row-wise to the single subquery result row.

See [Section 9.25.5](#) for details about the meaning of a row constructor comparison.

9.25. Row and Array Comparisons

This section describes several specialized constructs for making multiple comparisons between groups of values. These forms are syntactically related to the subquery forms of the previous section, but do not involve subqueries. The forms involving array subexpressions are Postgres Pro extensions; the rest are SQL-compliant. All of the expression forms documented in this section return Boolean (true/false) results.

9.25.1. IN

```
expression IN (value [, ...])
```

The right-hand side is a parenthesized list of expressions. The result is “true” if the left-hand expression's result is equal to any of the right-hand expressions. This is a shorthand notation for

```
expression = value1
OR
expression = value2
OR
...
```

Note that if the left-hand expression yields null, or if there are no equal right-hand values and at least one right-hand expression yields null, the result of the `IN` construct will be null, not false. This is in accordance with SQL's normal rules for Boolean combinations of null values.

9.25.2. NOT IN

```
expression NOT IN (value [, ...])
```

The right-hand side is a parenthesized list of expressions. The result is “true” if the left-hand expression's result is unequal to all of the right-hand expressions. This is a shorthand notation for

```
expression <> value1
AND
expression <> value2
AND
...
```

Note that if the left-hand expression yields null, or if there are no equal right-hand values and at least one right-hand expression yields null, the result of the `NOT IN` construct will be null, not true as one might naively expect. This is in accordance with SQL's normal rules for Boolean combinations of null values.

Tip

`x NOT IN y` is equivalent to `NOT (x IN y)` in all cases. However, null values are much more likely to trip up the novice when working with `NOT IN` than when working with `IN`. It is best to express your condition positively if possible.

9.25.3. ANY/SOME (array)

```
expression operator ANY (array expression)
expression operator SOME (array expression)
```

The right-hand side is a parenthesized expression, which must yield an array value. The left-hand expression is evaluated and compared to each element of the array using the given *operator*, which must yield a Boolean result. The result of `ANY` is “true” if any true result is obtained. The result is “false” if no true result is found (including the case where the array has zero elements).

If the array expression yields a null array, the result of `ANY` will be null. If the left-hand expression yields null, the result of `ANY` is ordinarily null (though a non-strict comparison operator could possibly yield a different result). Also, if the right-hand array contains any null elements and no true comparison result is obtained, the result of `ANY` will be null, not false (again, assuming a strict comparison operator). This is in accordance with SQL's normal rules for Boolean combinations of null values.

`SOME` is a synonym for `ANY`.

9.25.4. ALL (array)

```
expression operator ALL (array expression)
```

The right-hand side is a parenthesized expression, which must yield an array value. The left-hand expression is evaluated and compared to each element of the array using the given *operator*, which must yield a Boolean result. The result of `ALL` is “true” if all comparisons yield true (including the case where the array has zero elements). The result is “false” if any false result is found.

If the array expression yields a null array, the result of `ALL` will be null. If the left-hand expression yields null, the result of `ALL` is ordinarily null (though a non-strict comparison operator could possibly yield a different result). Also, if the right-hand array contains any null elements and no false comparison result is obtained, the result of `ALL` will be null, not true (again, assuming a strict comparison operator). This is in accordance with SQL's normal rules for Boolean combinations of null values.

9.25.5. Row Constructor Comparison

```
row_constructor operator row_constructor
```

Each side is a row constructor, as described in [Section 4.2.13](#). The two row constructors must have the same number of fields. The given *operator* is applied to each pair of corresponding fields. (Since the fields could be of different types, this means that a different specific operator could be selected for each pair.) All the selected operators must be members of some B-tree operator class, or be the negator of an = member of a B-tree operator class, meaning that row constructor comparison is only possible when the *operator* is `=`, `<>`, `<`, `<=`, `>`, or `>=`, or has semantics similar to one of these.

The `=` and `<>` cases work slightly differently from the others. Two rows are considered equal if all their corresponding members are non-null and equal; the rows are unequal if any corresponding members are non-null and unequal; otherwise the result of the row comparison is unknown (null).

For the `<`, `<=`, `>` and `>=` cases, the row elements are compared left-to-right, stopping as soon as an unequal or null pair of elements is found. If either of this pair of elements is null, the result of the row comparison is unknown (null); otherwise comparison of this pair of elements determines the result. For example, `ROW(1,2,NULL) < ROW(1,3,0)` yields true, not null, because the third pair of elements are not considered.

```
row_constructor IS DISTINCT FROM row_constructor
```

This construct is similar to a `<>` row comparison, but it does not yield null for null inputs. Instead, any null value is considered unequal to (distinct from) any non-null value, and any two nulls are considered equal (not distinct). Thus the result will either be true or false, never null.

```
row_constructor IS NOT DISTINCT FROM row_constructor
```

This construct is similar to a `=` row comparison, but it does not yield null for null inputs. Instead, any null value is considered unequal to (distinct from) any non-null value, and any two nulls are considered equal (not distinct). Thus the result will always be either true or false, never null.

9.25.6. Composite Type Comparison

```
record operator record
```

The SQL specification requires row-wise comparison to return NULL if the result depends on comparing two NULL values or a NULL and a non-NULL. Postgres Pro does this only when comparing the results of two row constructors (as in [Section 9.25.5](#)) or comparing a row constructor to the output of a subquery (as in [Section 9.24](#)). In other contexts where two composite-type values are compared, two NULL field values are considered equal, and a NULL is considered larger than a non-NULL. This is necessary in order to have consistent sorting and indexing behavior for composite types.

Each side is evaluated and they are compared row-wise. Composite type comparisons are allowed when the *operator* is `=`, `<>`, `<`, `<=`, `>` or `>=`, or has semantics similar to one of these. (To be specific, an operator can be a row comparison operator if it is a member of a B-tree operator class, or is the negator of the `=` member of a B-tree operator class.) The default behavior of the above operators is the same as for `IS [NOT] DISTINCT FROM` for row constructors (see [Section 9.25.5](#)).

To support matching of rows which include elements without a default B-tree operator class, the following operators are defined for composite type comparison: `*=`, `*<>`, `*<`, `*<=`, `*>`, and `*>=`. These operators compare the internal binary representation of the two rows. Two rows might have a different binary representation even though comparisons of the two rows with the equality operator is true. The ordering of rows under these comparison operators is deterministic but not otherwise meaningful. These operators are used internally for materialized views and might be useful for other specialized purposes such as replication and B-Tree deduplication (see [Section 65.1.4.3](#)). They are not intended to be generally useful for writing queries, though.

9.26. Set Returning Functions

This section describes functions that possibly return more than one row. The most widely used functions in this class are series generating functions, as detailed in [Table 9.67](#) and [Table 9.68](#). Other, more specialized set-returning functions are described elsewhere in this manual. See [Section 7.2.1.4](#) for ways to combine multiple set-returning functions.

Table 9.67. Series Generating Functions

Function	Description
<code>generate_series (start integer, stop integer [, step integer])</code>	<code>→ setof integer</code>
<code>generate_series (start bigint, stop bigint [, step bigint])</code>	<code>→ setof bigint</code>
<code>generate_series (start numeric, stop numeric [, step numeric])</code>	<code>→ setof numeric</code>
Generates a series of values from <i>start</i> to <i>stop</i> , with a step size of <i>step</i> . <i>step</i> defaults to 1.	
<code>generate_series (start timestamp, stop timestamp, step interval)</code>	<code>→ setof timestamp</code>

Function	Description
<code>generate_series (start timestamp with time zone, stop timestamp with time zone, step interval [, timezone text])</code>	<code>→ set of timestamp with time zone</code> Generates a series of values from <i>start</i> to <i>stop</i> , with a step size of <i>step</i> . In the time-zone-aware form, times of day and daylight-savings adjustments are computed according to the time zone named by the <i>timezone</i> argument, or the current TimeZone setting if that is omitted.

When *step* is positive, zero rows are returned if *start* is greater than *stop*. Conversely, when *step* is negative, zero rows are returned if *start* is less than *stop*. Zero rows are also returned if any input is NULL. It is an error for *step* to be zero. Some examples follow:

```
SELECT * FROM generate_series(2,4);
```

```
generate_series
-----
                2
                3
                4
(3 rows)
```

```
SELECT * FROM generate_series(5,1,-2);
```

```
generate_series
-----
                5
                3
                1
(3 rows)
```

```
SELECT * FROM generate_series(4,3);
```

```
generate_series
-----
(0 rows)
```

```
SELECT generate_series(1.1, 4, 1.3);
```

```
generate_series
-----
            1.1
            2.4
            3.7
(3 rows)
```

-- this example relies on the date-plus-integer operator:

```
SELECT current_date + s.a AS dates FROM generate_series(0,14,7) AS s(a);
```

```
dates
-----
2004-02-05
2004-02-12
2004-02-19
(3 rows)
```

```
SELECT * FROM generate_series('2008-03-01 00:00'::timestamp,
                              '2008-03-04 12:00', '10 hours');
```

```
generate_series
-----
2008-03-01 00:00:00
2008-03-01 10:00:00
2008-03-01 20:00:00
```

```

2008-03-02 06:00:00
2008-03-02 16:00:00
2008-03-03 02:00:00
2008-03-03 12:00:00
2008-03-03 22:00:00
2008-03-04 08:00:00
(9 rows)

-- this example assumes that TimeZone is set to UTC; note the DST transition:
SELECT * FROM generate_series('2001-10-22 00:00 -04:00'::timestampz,
                             '2001-11-01 00:00 -05:00'::timestampz,
                             '1 day'::interval, 'America/New_York');

    generate_series
-----
2001-10-22 04:00:00+00
2001-10-23 04:00:00+00
2001-10-24 04:00:00+00
2001-10-25 04:00:00+00
2001-10-26 04:00:00+00
2001-10-27 04:00:00+00
2001-10-28 04:00:00+00
2001-10-29 05:00:00+00
2001-10-30 05:00:00+00
2001-10-31 05:00:00+00
2001-11-01 05:00:00+00
(11 rows)

```

Table 9.68. Subscript Generating Functions

Function	Description
<code>generate_subscripts (array anyarray, dim integer) → setof integer</code>	Generates a series comprising the valid subscripts of the <i>dim</i> 'th dimension of the given array.
<code>generate_subscripts (array anyarray, dim integer, reverse boolean) → setof integer</code>	Generates a series comprising the valid subscripts of the <i>dim</i> 'th dimension of the given array. When <i>reverse</i> is true, returns the series in reverse order.

`generate_subscripts` is a convenience function that generates the set of valid subscripts for the specified dimension of the given array. Zero rows are returned for arrays that do not have the requested dimension, or if any input is `NULL`. Some examples follow:

```

-- basic usage:
SELECT generate_subscripts('{NULL,1,NULL,2}'::int[], 1) AS s;
 s
---
 1
 2
 3
 4
(4 rows)

-- presenting an array, the subscript and the subscripted
-- value requires a subquery:
SELECT * FROM arrays;
    a
-----
{-1,-2}
{100,200,300}

```

(2 rows)

```
SELECT a AS array, s AS subscript, a[s] AS value
FROM (SELECT generate_subscripts(a, 1) AS s, a FROM arrays) foo;
```

array	subscript	value
{-1,-2}	1	-1
{-1,-2}	2	-2
{100,200,300}	1	100
{100,200,300}	2	200
{100,200,300}	3	300

(5 rows)

```
-- unnest a 2D array:
CREATE OR REPLACE FUNCTION unnest2(anyarray)
RETURNS SETOF anyelement AS $$
select $1[i][j]
    from generate_subscripts($1,1) g1(i),
         generate_subscripts($1,2) g2(j);
$$ LANGUAGE sql IMMUTABLE;
CREATE FUNCTION
SELECT * FROM unnest2(ARRAY[[1,2],[3,4]]);
unnest2
```

1
2
3
4

(4 rows)

When a function in the `FROM` clause is suffixed by `WITH ORDINALITY`, a `bigint` column is appended to the function's output column(s), which starts from 1 and increments by 1 for each row of the function's output. This is most useful in the case of set returning functions such as `unnest()`.

```
-- set returning function WITH ORDINALITY:
SELECT * FROM pg_ls_dir('.') WITH ORDINALITY AS t(ls,n);
```

ls	n
pg_serial	1
pg_twophase	2
postmaster.opts	3
pg_notify	4
postgresql.conf	5
pg_tblspc	6
logfile	7
base	8
postmaster.pid	9
pg_ident.conf	10
global	11
pg_xact	12
pg_snapshots	13
pg_multixact	14
PG_VERSION	15
pg_wal	16
pg_hba.conf	17
pg_stat_tmp	18
pg_subtrans	19

(19 rows)

9.27. System Information Functions and Operators

The functions described in this section are used to obtain various information about a Postgres Pro installation.

9.27.1. Session Information Functions

[Table 9.69](#) shows several functions that extract session and system information.

In addition to the functions listed in this section, there are a number of functions related to the statistics system that also provide system information. See [Section 27.2.28](#) for more information.

Table 9.69. Session Information Functions

Function	Description
<code>current_catalog</code> \rightarrow name <code>current_database</code> $() \rightarrow$ name	Returns the name of the current database. (Databases are called “catalogs” in the SQL standard, so <code>current_catalog</code> is the standard's spelling.)
<code>current_query</code> $() \rightarrow$ text	Returns the text of the currently executing query, as submitted by the client (which might contain more than one statement).
<code>current_role</code> \rightarrow name	This is equivalent to <code>current_user</code> .
<code>current_schema</code> \rightarrow name <code>current_schema</code> $() \rightarrow$ name	Returns the name of the schema that is first in the search path (or a null value if the search path is empty). This is the schema that will be used for any tables or other named objects that are created without specifying a target schema.
<code>current_schemas</code> $(include_implicit \text{ boolean}) \rightarrow$ name[]	Returns an array of the names of all schemas presently in the effective search path, in their priority order. (Items in the current search_path setting that do not correspond to existing, searchable schemas are omitted.) If the Boolean argument is <code>true</code> , then implicitly-searched system schemas such as <code>pg_catalog</code> are included in the result.
<code>current_user</code> \rightarrow name	Returns the user name of the current execution context.
<code>inet_client_addr</code> $() \rightarrow$ inet	Returns the IP address of the current client, or <code>NULL</code> if the current connection is via a Unix-domain socket.
<code>inet_client_port</code> $() \rightarrow$ integer	Returns the IP port number of the current client, or <code>NULL</code> if the current connection is via a Unix-domain socket.
<code>inet_server_addr</code> $() \rightarrow$ inet	Returns the IP address on which the server accepted the current connection, or <code>NULL</code> if the current connection is via a Unix-domain socket.
<code>inet_server_port</code> $() \rightarrow$ integer	Returns the IP port number on which the server accepted the current connection, or <code>NULL</code> if the current connection is via a Unix-domain socket.
<code>pg_backend_pid</code> $() \rightarrow$ integer	Returns the process ID of the server process attached to the current session.

Function	Description
<code>pg_blocking_pids (integer) → integer[]</code>	<p>Returns an array of the process ID(s) of the sessions that are blocking the server process with the specified process ID from acquiring a lock, or an empty array if there is no such server process or it is not blocked.</p> <p>One server process blocks another if it either holds a lock that conflicts with the blocked process's lock request (hard block), or is waiting for a lock that would conflict with the blocked process's lock request and is ahead of it in the wait queue (soft block). When using parallel queries the result always lists client-visible process IDs (that is, <code>pg_backend_pid</code> results) even if the actual lock is held or awaited by a child worker process. As a result of that, there may be duplicated PIDs in the result. Also note that when a prepared transaction holds a conflicting lock, it will be represented by a zero process ID.</p> <p>Frequent calls to this function could have some impact on database performance, because it needs exclusive access to the lock manager's shared state for a short time.</p>
<code>pg_conf_load_time () → timestamp with time zone</code>	<p>Returns the time when the server configuration files were last loaded. If the current session was alive at the time, this will be the time when the session itself re-read the configuration files (so the reading will vary a little in different sessions). Otherwise it is the time when the postmaster process re-read the configuration files.</p>
<code>pg_current_logfile ([text]) → text</code>	<p>Returns the path name of the log file currently in use by the logging collector. The path includes the log_directory directory and the individual log file name. The result is <code>NULL</code> if the logging collector is disabled. When multiple log files exist, each in a different format, <code>pg_current_logfile</code> without an argument returns the path of the file having the first format found in the ordered list: <code>stderr</code>, <code>csvlog</code>, <code>jsonlog</code>. <code>NULL</code> is returned if no log file has any of these formats. To request information about a specific log file format, supply either <code>csvlog</code>, <code>jsonlog</code> or <code>stderr</code> as the value of the optional parameter. The result is <code>NULL</code> if the log format requested is not configured in log_destination. The result reflects the contents of the <code>current_logfiles</code> file.</p> <p>This function is restricted to superusers and roles with privileges of the <code>pg_monitor</code> role by default, but other users can be granted <code>EXECUTE</code> to run the function.</p>
<code>pg_my_temp_schema () → oid</code>	<p>Returns the OID of the current session's temporary schema, or zero if it has none (because it has not created any temporary tables).</p>
<code>pg_is_other_temp_schema (oid) → boolean</code>	<p>Returns true if the given OID is the OID of another session's temporary schema. (This can be useful, for example, to exclude other sessions' temporary tables from a catalog display.)</p>
<code>pg_jit_available () → boolean</code>	<p>Returns true if a JIT compiler extension is available (see Chapter 30) and the <code>jit</code> configuration parameter is set to <code>on</code>.</p>
<code>pg_listening_channels () → setof text</code>	<p>Returns the set of names of asynchronous notification channels that the current session is listening to.</p>
<code>pg_notification_queue_usage () → double precision</code>	<p>Returns the fraction (0-1) of the asynchronous notification queue's maximum size that is currently occupied by notifications that are waiting to be processed. See LISTEN and NOTIFY for more information.</p>
<code>pg_postmaster_start_time () → timestamp with time zone</code>	<p>Returns the time when the server started.</p>
<code>pg_safe_snapshot_blocking_pids (integer) → integer[]</code>	

Function	Description
	<p>Returns an array of the process ID(s) of the sessions that are blocking the server process with the specified process ID from acquiring a safe snapshot, or an empty array if there is no such server process or it is not blocked.</p> <p>A session running a <code>SERIALIZABLE</code> transaction blocks a <code>SERIALIZABLE READ ONLY DEFERRABLE</code> transaction from acquiring a snapshot until the latter determines that it is safe to avoid taking any predicate locks. See Section 13.2.3 for more information about serializable and deferrable transactions.</p> <p>Frequent calls to this function could have some impact on database performance, because it needs access to the predicate lock manager's shared state for a short time.</p>
<code>pg_trigger_depth</code> $() \rightarrow \text{integer}$	Returns the current nesting level of Postgres Pro triggers (0 if not called, directly or indirectly, from inside a trigger).
<code>session_user</code> $\rightarrow \text{name}$	Returns the session user's name.
<code>system_user</code> $\rightarrow \text{text}$	Returns the authentication method and the identity (if any) that the user presented during the authentication cycle before they were assigned a database role. It is represented as <code>auth_method:identity</code> or <code>NULL</code> if the user has not been authenticated (for example if Trust authentication has been used).
<code>user</code> $\rightarrow \text{name}$	This is equivalent to <code>current_user</code> .

Note

`current_catalog`, `current_role`, `current_schema`, `current_user`, `session_user`, and `user` have special syntactic status in SQL: they must be called without trailing parentheses. In Postgres Pro, parentheses can optionally be used with `current_schema`, but not with the others.

The `session_user` is normally the user who initiated the current database connection; but superusers can change this setting with [SET SESSION AUTHORIZATION](#). The `current_user` is the user identifier that is applicable for permission checking. Normally it is equal to the session user, but it can be changed with [SET ROLE](#). It also changes during the execution of functions with the attribute `SECURITY DEFINER`. In Unix parlance, the session user is the “real user” and the current user is the “effective user”. `current_role` and `user` are synonyms for `current_user`. (The SQL standard draws a distinction between `current_role` and `current_user`, but Postgres Pro does not, since it unifies users and roles into a single kind of entity.)

9.27.2. Access Privilege Inquiry Functions

[Table 9.70](#) lists functions that allow querying object access privileges programmatically. (See [Section 5.8](#) for more information about privileges.) In these functions, the user whose privileges are being inquired about can be specified by name or by OID (`pg_authid.oid`), or if the name is given as `public` then the privileges of the `PUBLIC` pseudo-role are checked. Also, the `user` argument can be omitted entirely, in which case the `current_user` is assumed. The object that is being inquired about can be specified either by name or by OID, too. When specifying by name, a schema name can be included if relevant. The access privilege of interest is specified by a text string, which must evaluate to one of the appropriate privilege keywords for the object's type (e.g., `SELECT`). Optionally, `WITH GRANT OPTION` can be added to a privilege type to test whether the privilege is held with grant option. Also, multiple privilege types can be listed separated by commas, in which case the result will be true if any of the listed privileges is held. (Case of the privilege string is not significant, and extra whitespace is allowed between but not within privilege names.) Some examples:

```
SELECT has_table_privilege('myschema.mytable', 'select');
SELECT has_table_privilege('joe', 'mytable', 'INSERT, SELECT WITH GRANT OPTION');
```

Table 9.70. Access Privilege Inquiry Functions

Function	Description
<code>has_any_column_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>table text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for any column of table? This succeeds either if the privilege is held for the whole table, or if there is a column-level grant of the privilege for at least one column. Allowable privilege types are SELECT, INSERT, UPDATE, and REFERENCES.
<code>has_column_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>table text</i> or <i>oid</i> , <i>column text</i> or <i>smallint</i> , <i>privilege text</i>) → boolean	Does user have privilege for the specified table column? This succeeds either if the privilege is held for the whole table, or if there is a column-level grant of the privilege for the column. The column can be specified by name or by attribute number (<code>pg_attribute.attnum</code>). Allowable privilege types are SELECT, INSERT, UPDATE, and REFERENCES.
<code>has_database_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>database text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for database? Allowable privilege types are CREATE, CONNECT, TEMPORARY, and TEMP (which is equivalent to TEMPORARY).
<code>has_foreign_data_wrapper_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>fdw text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for foreign-data wrapper? The only allowable privilege type is USAGE.
<code>has_function_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>function text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for function? The only allowable privilege type is EXECUTE. When specifying a function by name rather than by OID, the allowed input is the same as for the regprocedure data type (see Section 8.19). An example is: <code>SELECT has_function_privilege('joeuser', 'myfunc(int, text)', 'execute');</code>
<code>has_language_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>language text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for language? The only allowable privilege type is USAGE.
<code>has_parameter_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>parameter text</i> , <i>privilege text</i>) → boolean	Does user have privilege for configuration parameter? The parameter name is case-insensitive. Allowable privilege types are SET and ALTER SYSTEM.
<code>has_schema_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>schema text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for schema? Allowable privilege types are CREATE and USAGE.
<code>has_sequence_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>sequence text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for sequence? Allowable privilege types are USAGE, SELECT, and UPDATE.
<code>has_server_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>server text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for foreign server? The only allowable privilege type is USAGE.
<code>has_table_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>table text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	Does user have privilege for table? Allowable privilege types are SELECT, INSERT, UPDATE, DELETE, TRUNCATE, REFERENCES, TRIGGER, and MAINTAIN.
<code>has_tablespace_privilege</code> ([<i>user name</i> or <i>oid</i> ,] <i>tablespace text</i> or <i>oid</i> , <i>privilege text</i>) → boolean	

Function	Description
	Does user have privilege for tablespace? The only allowable privilege type is CREATE.
<code>has_type_privilege ([user name or oid,] type text or oid, privilege text) → boolean</code>	Does user have privilege for data type? The only allowable privilege type is USAGE. When specifying a type by name rather than by OID, the allowed input is the same as for the <code>reg-type</code> data type (see Section 8.19).
<code>pg_has_role ([user name or oid,] role text or oid, privilege text) → boolean</code>	Does user have privilege for role? Allowable privilege types are MEMBER, USAGE, and SET. MEMBER denotes direct or indirect membership in the role without regard to what specific privileges may be conferred. USAGE denotes whether the privileges of the role are immediately available without doing SET ROLE, while SET denotes whether it is possible to change to the role using the SET ROLE command. WITH ADMIN OPTION or WITH GRANT OPTION can be added to any of these privilege types to test whether the ADMIN privilege is held (all six spellings test the same thing). This function does not allow the special case of setting <code>user</code> to <code>public</code> , because the PUBLIC pseudo-role can never be a member of real roles.
<code>row_security_active (table text or oid) → boolean</code>	Is row-level security active for the specified table in the context of the current user and current environment?

[Table 9.71](#) shows the operators available for the `aclitem` type, which is the catalog representation of access privileges. See [Section 5.8](#) for information about how to read access privilege values.

Table 9.71. `aclitem` Operators

Operator	Description
<code>aclitem = aclitem → boolean</code>	Are <code>aclitems</code> equal? (Notice that type <code>aclitem</code> lacks the usual set of comparison operators; it has only equality. In turn, <code>aclitem</code> arrays can only be compared for equality.) 'calvin=r*w/hobbes'::aclitem = 'calvin=r*w*/hobbes'::aclitem → f
<code>aclitem[] @> aclitem → boolean</code>	Does array contain the specified privileges? (This is true if there is an array entry that matches the <code>aclitem</code> 's grantee and grantor, and has at least the specified set of privileges.) '{calvin=r*w/hobbes,hobbes=r*w*/postgres}'::aclitem[] @> 'calvin=r*/hobbes'::aclitem → t
<code>aclitem[] ~ aclitem → boolean</code>	This is a deprecated alias for <code>@></code> . '{calvin=r*w/hobbes,hobbes=r*w*/postgres}'::aclitem[] ~ 'calvin=r*/hobbes'::aclitem → t

[Table 9.72](#) shows some additional functions to manage the `aclitem` type.

Table 9.72. `aclitem` Functions

Function	Description
<code>acldefault (type "char", ownerId oid) → aclitem[]</code>	Constructs an <code>aclitem</code> array holding the default access privileges for an object of type <code>type</code> belonging to the role with OID <code>ownerId</code> . This represents the access privileges that will be assumed when an object's ACL entry is null. (The default access privileges are described in Section 5.8 .) The <code>type</code> parameter must be one of 'c' for COLUMN, 'r' for TABLE and table-like objects, 's' for SEQUENCE, 'd' for DATABASE, 'f' for FUNCTION or PROCEDURE, 'l' for LANGUAGE, 'L'

Function	Description
	for LARGE OBJECT, 'n' for SCHEMA, 'p' for PARAMETER, 't' for TABLESPACE, 'F' for FOREIGN DATA WRAPPER, 'S' for FOREIGN SERVER, or 'T' for TYPE or DOMAIN.
<code>aclxplode (aclitem[]) → setof record (grantor oid, grantee oid, privilege_type text, is_grantable boolean)</code>	Returns the <code>aclitem</code> array as a set of rows. If the grantee is the pseudo-role PUBLIC, it is represented by zero in the <code>grantee</code> column. Each granted privilege is represented as SELECT, INSERT, etc (see Table 5.1 for a full list). Note that each privilege is broken out as a separate row, so only one keyword appears in the <code>privilege_type</code> column.
<code>makeaclitem (grantee oid, grantor oid, privileges text, is_grantable boolean) → aclitem</code>	Constructs an <code>aclitem</code> with the given properties. <code>privileges</code> is a comma-separated list of privilege names such as SELECT, INSERT, etc, all of which are set in the result. (Case of the privilege string is not significant, and extra whitespace is allowed between but not within privilege names.)

9.27.3. Schema Visibility Inquiry Functions

[Table 9.73](#) shows functions that determine whether a certain object is *visible* in the current schema search path. For example, a table is said to be visible if its containing schema is in the search path and no table of the same name appears earlier in the search path. This is equivalent to the statement that the table can be referenced by name without explicit schema qualification. Thus, to list the names of all visible tables:

```
SELECT relname FROM pg_class WHERE pg_table_is_visible(oid);
```

For functions and operators, an object in the search path is said to be visible if there is no object of the same name *and argument data type(s)* earlier in the path. For operator classes and families, both the name and the associated index access method are considered.

Table 9.73. Schema Visibility Inquiry Functions

Function	Description
<code>pg_collation_is_visible (collation oid) → boolean</code> Is collation visible in search path?	
<code>pg_conversion_is_visible (conversion oid) → boolean</code> Is conversion visible in search path?	
<code>pg_function_is_visible (function oid) → boolean</code> Is function visible in search path? (This also works for procedures and aggregates.)	
<code>pg_opclass_is_visible (opclass oid) → boolean</code> Is operator class visible in search path?	
<code>pg_operator_is_visible (operator oid) → boolean</code> Is operator visible in search path?	
<code>pg_opfamily_is_visible (opclass oid) → boolean</code> Is operator family visible in search path?	
<code>pg_statistics_obj_is_visible (stat oid) → boolean</code> Is statistics object visible in search path?	
<code>pg_table_is_visible (table oid) → boolean</code> Is table visible in search path? (This works for all types of relations, including views, materialized views, indexes, sequences and foreign tables.)	
<code>pg_ts_config_is_visible (config oid) → boolean</code>	

Function	Description
	Is text search configuration visible in search path?
<code>pg_ts_dict_is_visible (dict oid) → boolean</code>	Is text search dictionary visible in search path?
<code>pg_ts_parser_is_visible (parser oid) → boolean</code>	Is text search parser visible in search path?
<code>pg_ts_template_is_visible (template oid) → boolean</code>	Is text search template visible in search path?
<code>pg_type_is_visible (type oid) → boolean</code>	Is type (or domain) visible in search path?

All these functions require object OIDs to identify the object to be checked. If you want to test an object by name, it is convenient to use the OID alias types (`regclass`, `regtype`, `regprocedure`, `regoperator`, `regconfig`, or `regdictionary`), for example:

```
SELECT pg_type_is_visible('myschema.widget'::regtype);
```

Note that it would not make much sense to test a non-schema-qualified type name in this way — if the name can be recognized at all, it must be visible.

9.27.4. System Catalog Information Functions

Table 9.74 lists functions that extract information from the system catalogs.

Table 9.74. System Catalog Information Functions

Function	Description
<code>format_type (type oid, typemod integer) → text</code>	Returns the SQL name for a data type that is identified by its type OID and possibly a type modifier. Pass NULL for the type modifier if no specific modifier is known.
<code>pg_basetype (regtype) → regtype</code>	Returns the OID of the base type of a domain identified by its type OID. If the argument is the OID of a non-domain type, returns the argument as-is. Returns NULL if the argument is not a valid type OID. If there's a chain of domain dependencies, it will recurse until finding the base type. Assuming <code>CREATE DOMAIN mytext AS text</code> : <code>pg_basetype('mytext'::regtype) → text</code>
<code>pg_char_to_encoding (encoding name) → integer</code>	Converts the supplied encoding name into an integer representing the internal identifier used in some system catalog tables. Returns -1 if an unknown encoding name is provided.
<code>pg_encoding_to_char (encoding integer) → name</code>	Converts the integer used as the internal identifier of an encoding in some system catalog tables into a human-readable string. Returns an empty string if an invalid encoding number is provided.
<code>pg_get_catalog_foreign_keys () → setof record (fktable regclass, fkcols text[], pktable regclass, pkcols text[], is_array boolean, is_opt boolean)</code>	Returns a set of records describing the foreign key relationships that exist within the Postgres Pro system catalogs. The <code>fktable</code> column contains the name of the referencing catalog, and the <code>fkcols</code> column contains the name(s) of the referencing column(s). Similarly, the <code>pktable</code> column contains the name of the referenced catalog, and the <code>pkcols</code> column contains the name(s) of the referenced column(s). If <code>is_array</code> is true, the last referencing column is an array, each of whose elements should match some entry in the referenced catalog. If <code>is_</code>

Function	Description
	<i>opt</i> is true, the referencing column(s) are allowed to contain zeroes instead of a valid reference.
<code>pg_get_constraintdef (<i>constraint oid</i> [, <i>pretty boolean</i>]) → text</code>	Reconstructs the creating command for a constraint. (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_expr (<i>expr pg_node_tree</i> , <i>relation oid</i> [, <i>pretty boolean</i>]) → text</code>	Decompiles the internal form of an expression stored in the system catalogs, such as the default value for a column. If the expression might contain Vars, specify the OID of the relation they refer to as the second parameter; if no Vars are expected, passing zero is sufficient.
<code>pg_get_functiondef (<i>func oid</i>) → text</code>	Reconstructs the creating command for a function or procedure. (This is a decompiled reconstruction, not the original text of the command.) The result is a complete <code>CREATE OR REPLACE FUNCTION</code> or <code>CREATE OR REPLACE PROCEDURE</code> statement.
<code>pg_get_function_arguments (<i>func oid</i>) → text</code>	Reconstructs the argument list of a function or procedure, in the form it would need to appear in within <code>CREATE FUNCTION</code> (including default values).
<code>pg_get_function_identity_arguments (<i>func oid</i>) → text</code>	Reconstructs the argument list necessary to identify a function or procedure, in the form it would need to appear in within commands such as <code>ALTER FUNCTION</code> . This form omits default values.
<code>pg_get_function_result (<i>func oid</i>) → text</code>	Reconstructs the <code>RETURNS</code> clause of a function, in the form it would need to appear in within <code>CREATE FUNCTION</code> . Returns <code>NULL</code> for a procedure.
<code>pg_get_indexdef (<i>index oid</i> [, <i>column integer</i>, <i>pretty boolean</i>]) → text</code>	Reconstructs the creating command for an index. (This is a decompiled reconstruction, not the original text of the command.) If <i>column</i> is supplied and is not zero, only the definition of that column is reconstructed.
<code>pg_get_keywords () → setof record (<i>word text</i>, <i>catcode "char"</i>, <i>barelabel boolean</i>, <i>catdesc text</i>, <i>baredesc text</i>)</code>	Returns a set of records describing the SQL keywords recognized by the server. The <i>word</i> column contains the keyword. The <i>catcode</i> column contains a category code: <code>U</code> for an unreserved keyword, <code>C</code> for a keyword that can be a column name, <code>T</code> for a keyword that can be a type or function name, or <code>R</code> for a fully reserved keyword. The <i>barelabel</i> column contains <code>true</code> if the keyword can be used as a “bare” column label in <code>SELECT</code> lists, or <code>false</code> if it can only be used after <code>AS</code> . The <i>catdesc</i> column contains a possibly-localized string describing the keyword's category. The <i>baredesc</i> column contains a possibly-localized string describing the keyword's column label status.
<code>pg_get_partkeydef (<i>table oid</i>) → text</code>	Reconstructs the definition of a partitioned table's partition key, in the form it would have in the <code>PARTITION BY</code> clause of <code>CREATE TABLE</code> . (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_ruledef (<i>rule oid</i> [, <i>pretty boolean</i>]) → text</code>	Reconstructs the creating command for a rule. (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_serial_sequence (<i>table text</i>, <i>column text</i>) → text</code>	Returns the name of the sequence associated with a column, or <code>NULL</code> if no sequence is associated with the column. If the column is an identity column, the associated sequence is the sequence internally created for that column. For columns created using one of the serial types

Function	Description
	<p>(<i>serial</i>, <i>smallserial</i>, <i>bigserial</i>), it is the sequence created for that serial column definition. In the latter case, the association can be modified or removed with <code>ALTER SEQUENCE OWNED BY</code>. (This function probably should have been called <code>pg_get_owned_sequence</code> ; its current name reflects the fact that it has historically been used with serial-type columns.) The first parameter is a table name with optional schema, and the second parameter is a column name. Because the first parameter potentially contains both schema and table names, it is parsed per usual SQL rules, meaning it is lower-cased by default. The second parameter, being just a column name, is treated literally and so has its case preserved. The result is suitably formatted for passing to the sequence functions (see Section 9.17).</p> <p>A typical use is in reading the current value of the sequence for an identity or serial column, for example:</p> <pre>SELECT currval(pg_get_serial_sequence('sometable', 'id'));</pre>
<code>pg_get_statisticsobjdef (<i>statobj oid</i>) → text</code>	Reconstructs the creating command for an extended statistics object. (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_triggerdef (<i>trigger oid</i> [, <i>pretty boolean</i>]) → text</code>	Reconstructs the creating command for a trigger. (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_userbyid (<i>role oid</i>) → name</code>	Returns a role's name given its OID.
<code>pg_get_viewdef (<i>view oid</i> [, <i>pretty boolean</i>]) → text</code>	Reconstructs the underlying <code>SELECT</code> command for a view or materialized view. (This is a decompiled reconstruction, not the original text of the command.)
<code>pg_get_viewdef (<i>view oid</i>, <i>wrap_column integer</i>) → text</code>	Reconstructs the underlying <code>SELECT</code> command for a view or materialized view. (This is a decompiled reconstruction, not the original text of the command.) In this form of the function, pretty-printing is always enabled, and long lines are wrapped to try to keep them shorter than the specified number of columns.
<code>pg_get_viewdef (<i>view text</i> [, <i>pretty boolean</i>]) → text</code>	Reconstructs the underlying <code>SELECT</code> command for a view or materialized view, working from a textual name for the view rather than its OID. (This is deprecated; use the OID variant instead.)
<code>pg_index_column_has_property (<i>index regclass</i>, <i>column integer</i>, <i>property text</i>) → boolean</code>	Tests whether an index column has the named property. Common index column properties are listed in Table 9.75 . (Note that extension access methods can define additional property names for their indexes.) <code>NULL</code> is returned if the property name is not known or does not apply to the particular object, or if the OID or column number does not identify a valid object.
<code>pg_index_has_property (<i>index regclass</i>, <i>property text</i>) → boolean</code>	Tests whether an index has the named property. Common index properties are listed in Table 9.76 . (Note that extension access methods can define additional property names for their indexes.) <code>NULL</code> is returned if the property name is not known or does not apply to the particular object, or if the OID does not identify a valid object.
<code>pg_indexam_has_property (<i>am oid</i>, <i>property text</i>) → boolean</code>	Tests whether an index access method has the named property. Access method properties are listed in Table 9.77 . <code>NULL</code> is returned if the property name is not known or does not apply to the particular object, or if the OID does not identify a valid object.

Function	Description
<code>pg_options_to_table (options_array text[]) → setof record (option_name text, option_value text)</code>	Returns the set of storage options represented by a value from <code>pg_class .reloptions</code> or <code>pg_attribute .attoptions</code> .
<code>pg_settings_get_flags (guc text) → text[]</code>	Returns an array of the flags associated with the given GUC, or <code>NULL</code> if it does not exist. The result is an empty array if the GUC exists but there are no flags to show. Only the most useful flags listed in Table 9.78 are exposed.
<code>pg_tablespace_databases (tablespace oid) → setof oid</code>	Returns the set of OIDs of databases that have objects stored in the specified tablespace. If this function returns any rows, the tablespace is not empty and cannot be dropped. To identify the specific objects populating the tablespace, you will need to connect to the database(s) identified by <code>pg_tablespace_databases</code> and query their <code>pg_class</code> catalogs.
<code>pg_tablespace_location (tablespace oid) → text</code>	Returns the file system path that this tablespace is located in.
<code>pg_typeof ("any") → regtype</code> <code>pg_typeof(33) → integer</code>	Returns the OID of the data type of the value that is passed to it. This can be helpful for troubleshooting or dynamically constructing SQL queries. The function is declared as returning <code>regtype</code> , which is an OID alias type (see Section 8.19); this means that it is the same as an OID for comparison purposes but displays as a type name.
<code>COLLATION FOR ("any") → text</code> <code>collation for ('foo'::text) → "default"</code> <code>collation for ('foo' COLLATE "de_DE") → "de_DE"</code>	Returns the name of the collation of the value that is passed to it. The value is quoted and schema-qualified if necessary. If no collation was derived for the argument expression, then <code>NULL</code> is returned. If the argument is not of a collatable data type, then an error is raised.
<code>to_regclass (text) → regclass</code>	Translates a textual relation name to its OID. A similar result is obtained by casting the string to type <code>regclass</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regcollation (text) → regcollation</code>	Translates a textual collation name to its OID. A similar result is obtained by casting the string to type <code>regcollation</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regnamespace (text) → regnamespace</code>	Translates a textual schema name to its OID. A similar result is obtained by casting the string to type <code>regnamespace</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regoper (text) → regoper</code>	Translates a textual operator name to its OID. A similar result is obtained by casting the string to type <code>regoper</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found or is ambiguous.
<code>to_regoperator (text) → regoperator</code>	Translates a textual operator name (with parameter types) to its OID. A similar result is obtained by casting the string to type <code>regoperator</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.

Function	Description
<code>to_regproc (text) → regproc</code>	Translates a textual function or procedure name to its OID. A similar result is obtained by casting the string to type <code>regproc</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found or is ambiguous.
<code>to_regprocedure (text) → regprocedure</code>	Translates a textual function or procedure name (with argument types) to its OID. A similar result is obtained by casting the string to type <code>regprocedure</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regprofile (text) → regprofile</code>	Translates a textual profile name to its OID. A similar result is obtained by casting the string to type <code>regprofile</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regrole (text) → regrole</code>	Translates a textual role name to its OID. A similar result is obtained by casting the string to type <code>regrole</code> (see Section 8.19); however, this function will return <code>NULL</code> rather than throwing an error if the name is not found.
<code>to_regtype (text) → regtype</code>	Parses a string of text, extracts a potential type name from it, and translates that name into a type OID. A syntax error in the string will result in an error; but if the string is a syntactically valid type name that happens not to be found in the catalogs, the result is <code>NULL</code> . A similar result is obtained by casting the string to type <code>regtype</code> (see Section 8.19), except that that will throw error for name not found.
<code>to_regtypemod (text) → integer</code>	<p>Parses a string of text, extracts a potential type name from it, and translates its type modifier, if any. A syntax error in the string will result in an error; but if the string is a syntactically valid type name that happens not to be found in the catalogs, the result is <code>NULL</code>. The result is <code>-1</code> if no type modifier is present.</p> <p><code>to_regtypemod</code> can be combined with <code>to_regtype</code> to produce appropriate inputs for <code>format_type</code>, allowing a string representing a type name to be canonicalized.</p> <pre>format_type(to_regtype('varchar(32)'), to_regtypemod('varchar(32)')) → character_varying(32)</pre>

Most of the functions that reconstruct (decompile) database objects have an optional *pretty* flag, which if `true` causes the result to be “pretty-printed”. Pretty-printing suppresses unnecessary parentheses and adds whitespace for legibility. The pretty-printed format is more readable, but the default format is more likely to be interpreted the same way by future versions of Postgres Pro; so avoid using pretty-printed output for dump purposes. Passing `false` for the *pretty* parameter yields the same result as omitting the parameter.

Table 9.75. Index Column Properties

Name	Description
<code>asc</code>	Does the column sort in ascending order on a forward scan?
<code>desc</code>	Does the column sort in descending order on a forward scan?
<code>nulls_first</code>	Does the column sort with nulls first on a forward scan?
<code>nulls_last</code>	Does the column sort with nulls last on a forward scan?

Name	Description
orderable	Does the column possess any defined sort ordering?
distance_orderable	Can the column be scanned in order by a “distance” operator, for example <code>ORDER BY col <-> constant</code> ?
returnable	Can the column value be returned by an index-only scan?
search_array	Does the column natively support <code>col = ANY(array)</code> searches?
search_nulls	Does the column support <code>IS NULL</code> and <code>IS NOT NULL</code> searches?

Table 9.76. Index Properties

Name	Description
clusterable	Can the index be used in a <code>CLUSTER</code> command?
index_scan	Does the index support plain (non-bitmap) scans?
bitmap_scan	Does the index support bitmap scans?
backward_scan	Can the scan direction be changed in mid-scan (to support <code>FETCH BACKWARD</code> on a cursor without needing materialization)?

Table 9.77. Index Access Method Properties

Name	Description
can_order	Does the access method support <code>ASC</code> , <code>DESC</code> and related keywords in <code>CREATE INDEX</code> ?
can_unique	Does the access method support unique indexes?
can_multi_col	Does the access method support indexes with multiple columns?
can_exclude	Does the access method support exclusion constraints?
can_include	Does the access method support the <code>INCLUDE</code> clause of <code>CREATE INDEX</code> ?

Table 9.78. GUC Flags

Flag	Description
EXPLAIN	Parameters with this flag are included in <code>EXPLAIN (SETTINGS)</code> commands.
NO_SHOW_ALL	Parameters with this flag are excluded from <code>SHOW ALL</code> commands.
NO_RESET	Parameters with this flag do not support <code>RESET</code> commands.
NO_RESET_ALL	Parameters with this flag are excluded from <code>RESET ALL</code> commands.
NOT_IN_SAMPLE	Parameters with this flag are not included in <code>postgresql.conf</code> by default.
RUNTIME_COMPUTED	Parameters with this flag are runtime-computed ones.

9.27.5. Object Information and Addressing Functions

Table 9.79 lists functions related to database object identification and addressing.

Table 9.79. Object Information and Addressing Functions

Function	Description
<code>pg_describe_object (classid oid, objid oid, objsubid integer) → text</code>	Returns a textual description of a database object identified by catalog OID, object OID, and sub-object ID (such as a column number within a table; the sub-object ID is zero when referring to a whole object). This description is intended to be human-readable, and might be translated, depending on server configuration. This is especially useful to determine the identity of an object referenced in the <code>pg_depend</code> catalog. This function returns <code>NULL</code> values for undefined objects.
<code>pg_identify_object (classid oid, objid oid, objsubid integer) → record (type text, schema text, name text, identity text)</code>	Returns a row containing enough information to uniquely identify the database object specified by catalog OID, object OID and sub-object ID. This information is intended to be machine-readable, and is never translated. <i>type</i> identifies the type of database object; <i>schema</i> is the schema name that the object belongs in, or <code>NULL</code> for object types that do not belong to schemas; <i>name</i> is the name of the object, quoted if necessary, if the name (along with schema name, if pertinent) is sufficient to uniquely identify the object, otherwise <code>NULL</code> ; <i>identity</i> is the complete object identity, with the precise format depending on object type, and each name within the format being schema-qualified and quoted as necessary. Undefined objects are identified with <code>NULL</code> values.
<code>pg_identify_object_as_address (classid oid, objid oid, objsubid integer) → record (type text, object_names text[], object_args text[])</code>	Returns a row containing enough information to uniquely identify the database object specified by catalog OID, object OID and sub-object ID. The returned information is independent of the current server, that is, it could be used to identify an identically named object in another server. <i>type</i> identifies the type of database object; <i>object_names</i> and <i>object_args</i> are text arrays that together form a reference to the object. These three values can be passed to <code>pg_get_object_address</code> to obtain the internal address of the object.
<code>pg_get_object_address (type text, object_names text[], object_args text[]) → record (classid oid, objid oid, objsubid integer)</code>	Returns a row containing enough information to uniquely identify the database object specified by a type code and object name and argument arrays. The returned values are the ones that would be used in system catalogs such as <code>pg_depend</code> ; they can be passed to other system functions such as <code>pg_describe_object</code> or <code>pg_identify_object</code> . <i>classid</i> is the OID of the system catalog containing the object; <i>objid</i> is the OID of the object itself, and <i>objsubid</i> is the sub-object ID, or zero if none. This function is the inverse of <code>pg_identify_object_as_address</code> . Undefined objects are identified with <code>NULL</code> values.

9.27.6. Comment Information Functions

The functions shown in Table 9.80 extract comments previously stored with the `COMMENT` command. A null value is returned if no comment could be found for the specified parameters.

Table 9.80. Comment Information Functions

Function	Description
<code>col_description (table oid, column integer) → text</code>	

Function	Description
	Returns the comment for a table column, which is specified by the OID of its table and its column number. (<code>obj_description</code> cannot be used for table columns, since columns do not have OIDs of their own.)
<code>obj_description (object oid, catalog name) → text</code>	Returns the comment for a database object specified by its OID and the name of the containing system catalog. For example, <code>obj_description(123456, 'pg_class')</code> would retrieve the comment for the table with OID 123456.
<code>obj_description (object oid) → text</code>	Returns the comment for a database object specified by its OID alone. This is <i>deprecated</i> since there is no guarantee that OIDs are unique across different system catalogs; therefore, the wrong comment might be returned.
<code>shobj_description (object oid, catalog name) → text</code>	Returns the comment for a shared database object specified by its OID and the name of the containing system catalog. This is just like <code>obj_description</code> except that it is used for retrieving comments on shared objects (that is, databases, roles, and tablespaces). Some system catalogs are global to all databases within each cluster, and the descriptions for objects in them are stored globally as well.

9.27.7. Data Validity Checking Functions

The functions shown in [Table 9.81](#) can be helpful for checking validity of proposed input data.

Table 9.81. Data Validity Checking Functions

Function	Description	Example(s)
<code>pg_input_is_valid</code>	<p>(<i>string</i> text, <i>type</i> text) → boolean</p> <p>Tests whether the given <i>string</i> is valid input for the specified data type, returning true or false.</p> <p>This function will only work as desired if the data type's input function has been updated to report invalid input as a “soft” error. Otherwise, invalid input will abort the transaction, just as if the string had been cast to the type directly.</p>	<pre>pg_input_is_valid('42', 'integer') → t pg_input_is_valid('420000000000', 'integer') → f pg_input_is_valid('1234.567', 'numeric(7,4)') → f</pre>
<code>pg_input_error_info</code>	<p>(<i>string</i> text, <i>type</i> text) → record (<i>message</i> text, <i>detail</i> text, <i>hint</i> text, <i>sql_error_code</i> text)</p> <p>Tests whether the given <i>string</i> is valid input for the specified data type; if not, return the details of the error that would have been thrown. If the input is valid, the results are NULL. The inputs are the same as for <code>pg_input_is_valid</code>.</p> <p>This function will only work as desired if the data type's input function has been updated to report invalid input as a “soft” error. Otherwise, invalid input will abort the transaction, just as if the string had been cast to the type directly.</p>	<pre>SELECT * FROM pg_input_error_info('420000000000', 'integer') → message detail hint sql_error_code -----+-----+-----+ +-----+ value "420000000000" is out of range for type integer 22003</pre>

9.27.8. Transaction ID and Snapshot Information Functions

The functions shown in [Table 9.82](#) provide server transaction information in an exportable form. The main use of these functions is to determine which transactions were committed between two snapshots.

Table 9.82. Transaction ID and Snapshot Information Functions

Function	Description
<code>age (xid) → integer</code>	Returns the number of transactions between the supplied transaction id and the current transaction counter.
<code>mxid_age (xid) → integer</code>	Returns the number of multixacts IDs between the supplied multixact ID and the current multixacts counter.
<code>pg_current_xact_id () → xid8</code>	Returns the current transaction's ID. It will assign a new one if the current transaction does not have one already (because it has not performed any database updates); see Section 67.1 for details. If executed in a subtransaction, this will return the top-level transaction ID; see Section 67.3 for details.
<code>pg_current_xact_id_if_assigned () → xid8</code>	Returns the current transaction's ID, or <code>NULL</code> if no ID is assigned yet. (It's best to use this variant if the transaction might otherwise be read-only, to avoid unnecessary consumption of an XID.) If executed in a subtransaction, this will return the top-level transaction ID.
<code>pg_xact_status (xid8) → text</code>	Reports the commit status of a recent transaction. The result is one of <code>in progress</code> , <code>committed</code> , or <code>aborted</code> , provided that the transaction is recent enough that the system retains the commit status of that transaction. If it is old enough that no references to the transaction survive in the system and the commit status information has been discarded, the result is <code>NULL</code> . Applications might use this function, for example, to determine whether their transaction committed or aborted after the application and database server become disconnected while a <code>COMMIT</code> is in progress. Note that prepared transactions are reported as <code>in progress</code> ; applications must check <code>pg_prepared_xacts</code> if they need to determine whether a transaction ID belongs to a prepared transaction.
<code>pg_current_snapshot () → pg_snapshot</code>	Returns a current <i>snapshot</i> , a data structure showing which transaction IDs are now in-progress. Only top-level transaction IDs are included in the snapshot; subtransaction IDs are not shown; see Section 67.3 for details.
<code>pg_snapshot_xip (pg_snapshot) → setof xid8</code>	Returns the set of in-progress transaction IDs contained in a snapshot.
<code>pg_snapshot_xmax (pg_snapshot) → xid8</code>	Returns the <code>xmax</code> of a snapshot.
<code>pg_snapshot_xmin (pg_snapshot) → xid8</code>	Returns the <code>xmin</code> of a snapshot.
<code>pg_visible_in_snapshot (xid8, pg_snapshot) → boolean</code>	Is the given transaction ID <i>visible</i> according to this snapshot (that is, was it completed before the snapshot was taken)? Note that this function will not give the correct answer for a subtransaction ID (subxid); see Section 67.3 for details.

In Postgres Pro Shardman, transaction IDs are implemented as 64-bit counters to prevent transaction ID wraparound. The functions shown in [Table 9.82](#), except `age` and `mxid_age`, use a 64-bit type `xid8` that does not wrap around during the life of an installation and can be converted to `xid` by casting if

required; see [Section 67.1](#) for details. The data type `pg_snapshot` stores information about transaction ID visibility at a particular moment in time. Its components are described in [Table 9.83](#). `pg_snapshot`'s textual representation is `xmin:xmax:xip_list`. For example `10:20:10,14,15` means `xmin=10`, `xmax=20`, `xip_list=10, 14, 15`.

Table 9.83. Snapshot Components

Name	Description
<code>xmin</code>	Lowest transaction ID that was still active. All transaction IDs less than <code>xmin</code> are either committed and visible, or rolled back and dead.
<code>xmax</code>	One past the highest completed transaction ID. All transaction IDs greater than or equal to <code>xmax</code> had not yet completed as of the time of the snapshot, and thus are invisible.
<code>xip_list</code>	Transactions in progress at the time of the snapshot. A transaction ID that is <code>xmin <= X < xmax</code> and not in this list was already completed at the time of the snapshot, and thus is either visible or dead according to its commit status. This list does not include the transaction IDs of subtransactions (subxids).

In releases of Postgres Pro before 13 there was no `xid8` type, so variants of these functions were provided that used `bigint` to represent a 64-bit XID, with a correspondingly distinct snapshot data type `txid_snapshot`. These older functions have `txid` in their names. They are still supported for backward compatibility, but may be removed from a future release. See [Table 9.84](#).

Table 9.84. Deprecated Transaction ID and Snapshot Information Functions

Function	Description
<code>txid_current () → bigint</code> See <code>pg_current_xact_id()</code> .	
<code>txid_current_if_assigned () → bigint</code> See <code>pg_current_xact_id_if_assigned()</code> .	
<code>txid_current_snapshot () → txid_snapshot</code> See <code>pg_current_snapshot()</code> .	
<code>txid_snapshot_xip (txid_snapshot) → setof bigint</code> See <code>pg_snapshot_xip()</code> .	
<code>txid_snapshot_xmax (txid_snapshot) → bigint</code> See <code>pg_snapshot_xmax()</code> .	
<code>txid_snapshot_xmin (txid_snapshot) → bigint</code> See <code>pg_snapshot_xmin()</code> .	
<code>txid_visible_in_snapshot (bigint, txid_snapshot) → boolean</code> See <code>pg_visible_in_snapshot()</code> .	
<code>txid_status (bigint) → text</code> See <code>pg_xact_status()</code> .	

9.27.9. Committed Transaction Information Functions

The functions shown in [Table 9.85](#) provide information about when past transactions were committed. They only provide useful data when the `track_commit_timestamp` configuration option is enabled, and

only for transactions that were committed after it was enabled. Commit timestamp information is routinely removed during vacuum.

Table 9.85. Committed Transaction Information Functions

Function	Description
<code>pg_xact_commit_timestamp (xid) → timestamp with time zone</code> Returns the commit timestamp of a transaction.	
<code>pg_xact_commit_timestamp_origin (xid) → record (timestamp with time zone, roident oid)</code> Returns the commit timestamp and replication origin of a transaction.	
<code>pg_last_committed_xact () → record (xid, timestamp with time zone, roident oid)</code> Returns the transaction ID, commit timestamp and replication origin of the latest committed transaction.	

9.27.10. Control Data Functions

The functions shown in [Table 9.86](#) print information initialized during `initdb`, such as the catalog version. They also show information about write-ahead logging and checkpoint processing. This information is cluster-wide, not specific to any one database.

Table 9.86. Control Data Functions

Function	Description
<code>pg_control_checkpoint () → record</code> Returns information about current checkpoint state, as shown in Table 9.87 .	
<code>pg_control_system () → record</code> Returns information about current control file state, as shown in Table 9.88 .	
<code>pg_control_init () → record</code> Returns information about cluster initialization state, as shown in Table 9.89 .	
<code>pg_control_recovery () → record</code> Returns information about recovery state, as shown in Table 9.90 .	

Table 9.87. pg_control_checkpoint Output Columns

Column Name	Data Type
<code>checkpoint_lsn</code>	<code>pg_lsn</code>
<code>redo_lsn</code>	<code>pg_lsn</code>
<code>redo_wal_file</code>	<code>text</code>
<code>timeline_id</code>	<code>integer</code>
<code>prev_timeline_id</code>	<code>integer</code>
<code>full_page_writes</code>	<code>boolean</code>
<code>next_xid</code>	<code>text</code>
<code>next_oid</code>	<code>oid</code>
<code>next_multixact_id</code>	<code>xid</code>
<code>next_multi_offset</code>	<code>xid</code>
<code>oldest_xid</code>	<code>xid</code>
<code>oldest_xid_dbid</code>	<code>oid</code>

Column Name	Data Type
oldest_active_xid	xid
oldest_multi_xid	xid
oldest_multi_dbid	oid
oldest_commit_ts_xid	xid
newest_commit_ts_xid	xid
checkpoint_time	timestamp with time zone

Table 9.88. pg_control_system Output Columns

Column Name	Data Type
pg_control_version	integer
catalog_version_no	integer
system_identifier	bigint
pg_control_last_modified	timestamp with time zone
pg_control_edition	text

Table 9.89. pg_control_init Output Columns

Column Name	Data Type
max_data_alignment	integer
database_block_size	integer
blocks_per_segment	integer
wal_block_size	integer
bytes_per_wal_segment	integer
max_identifier_length	integer
max_index_columns	integer
max_toast_chunk_size	integer
large_object_chunk_size	integer
float8_pass_by_value	boolean
data_page_checksum_version	integer

Table 9.90. pg_control_recovery Output Columns

Column Name	Data Type
min_recovery_end_lsn	pg_lsn
min_recovery_end_timeline	integer
backup_start_lsn	pg_lsn
backup_end_lsn	pg_lsn
end_of_backup_record_required	boolean

9.27.11. Version Information Functions

The functions shown in [Table 9.91](#) print version information.

Table 9.91. Version Information Functions

Function Description
version () → text

Function	Description
	Returns a string describing the PostgreSQL server's version. You can also get this information from server_version , or for a machine-readable version use server_version_num . Software developers should use <code>server_version_num</code> (available since 8.2) or <code>PQserverVersion</code> instead of parsing the text version.
<code>unicode_version () → text</code>	Returns a string representing the version of Unicode used by Postgres Pro.
<code>icu_unicode_version () → text</code>	Returns a string representing the version of Unicode used by ICU, if the server was built with ICU support; otherwise returns <code>NULL</code>
<code>pgpro_version () → text</code>	Returns a string describing the Postgres Pro server's version. This function is deprecated and will be removed in future releases. Use the pgpro_version configuration parameter instead.
<code>pgpro_edition () → text</code>	Returns a string describing Postgres Pro edition, i.e. <code>standard</code> or <code>enterprise</code> . This function is deprecated and will be removed in future releases. Use the pgpro_edition configuration parameter instead.
<code>pgpro_build () → text</code>	Returns the commit ID of Postgres Pro source files. This function is deprecated and will be removed in future releases. Use the pgpro_build configuration parameter instead.

9.27.12. WAL Summarization Information Functions

The functions shown in [Table 9.92](#) print information about the status of WAL summarization. See [summarize_wal](#).

Table 9.92. WAL Summarization Information Functions

Function	Description
<code>pg_available_wal_summaries (tli bigint, start_lsn pg_lsn , end_lsn pg_lsn)</code>	Returns information about the WAL summary files present in the data directory, under <code>pg_wal/summaries</code> . One row will be returned per WAL summary file. Each file summarizes WAL on the indicated TLI within the indicated LSN range. This function might be useful to determine whether enough WAL summaries are present on the server to take an incremental backup based on some prior backup whose start LSN is known.
<code>pg_wal_summary_contents (tli bigint, start_lsn pg_lsn , end_lsn pg_lsn) → setof record (relfilenode oid, reltablespace oid, reldatabase oid, relforknumber smallint, relblocknumber bigint, is_limit_block boolean)</code>	Returns one information about the contents of a single WAL summary file identified by TLI and starting and ending LSNs. Each row with <code>is_limit_block</code> <code>false</code> indicates that the block identified by the remaining output columns was modified by at least one WAL record within the range of records summarized by this file. Each row with <code>is_limit_block</code> <code>true</code> indicates either that (a) the relation fork was truncated to the length given by <code>relblocknumber</code> within the relevant range of WAL records or (b) that the relation fork was created or dropped within the relevant range of WAL records; in such cases, <code>relblocknumber</code> will be zero.
<code>pg_get_wal_summarizer_state () → record (summarized_tli bigint, summarized_lsn pg_lsn , pending_lsn pg_lsn , summarizer_pid int)</code>	Returns information about the progress of the WAL summarizer. If the WAL summarizer has never run since the instance was started, then <code>summarized_tli</code> and <code>summarized_lsn</code> will be 0 and 0/0 respectively; otherwise, they will be the TLI and ending LSN of the last WAL

Function	Description
	<p>summary file written to disk. If the WAL summarizer is currently running, <code>pending_lsn</code> will be the ending LSN of the last record that it has consumed, which must always be greater than or equal to <code>summarized_lsn</code> ; if the WAL summarizer is not running, it will be equal to <code>summarized_lsn</code> . <code>summarizer_pid</code> is the PID of the WAL summarizer process, if it is running, and otherwise NULL.</p> <p>As a special exception, the WAL summarizer will refuse to generate WAL summary files if run on WAL generated under <code>wal_level=minimal</code> , since such summaries would be unsafe to use as the basis for an incremental backup. In this case, the fields above will continue to advance as if summaries were being generated, but nothing will be written to disk. Once the summarizer reaches WAL generated while <code>wal_level</code> was set to <code>replica</code> or higher, it will resume writing summaries to disk.</p>

9.28. System Administration Functions

The functions described in this section are used to control and monitor a Postgres Pro installation.

9.28.1. Configuration Settings Functions

[Table 9.93](#) shows the functions available to query and alter run-time configuration parameters.

Table 9.93. Configuration Settings Functions

Function	Description
Example(s)	
<code>current_setting (<i>setting_name</i> text [, <i>missing_ok</i> boolean]) → text</code>	Returns the current value of the setting <i>setting_name</i> . If there is no such setting, <code>current_setting</code> throws an error unless <i>missing_ok</i> is supplied and is <code>true</code> (in which case NULL is returned). This function corresponds to the SQL command SHOW .
<code>current_setting('datestyle') → ISO, MDY</code>	
<code>set_config (<i>setting_name</i> text, <i>new_value</i> text, <i>is_local</i> boolean) → text</code>	Sets the parameter <i>setting_name</i> to <i>new_value</i> , and returns that value. If <i>is_local</i> is <code>true</code> , the new value will only apply during the current transaction. If you want the new value to apply for the rest of the current session, use <code>false</code> instead. This function corresponds to the SQL command SET .
<code>set_config('log_statement_stats', 'off', false) → off</code>	
<code>pg_backend_set_config (<i>pid</i> int, <i>config</i> text, <i>wait</i> int default 0) → boolean</code>	Sets one or more parameters provided in the <i>config</i> string for the backend with the specified PID. Following the <code>postgresql.conf</code> convention, each configuration parameter in the <i>config</i> string must be written on a separate line, in the <i>name=value</i> format. The <code>pg_backend_set_config</code> function only affects run-time parameters. When called, it modifies the main configuration until the end of the session or the next <code>pg_backend_set_config</code> call. The changes take effect starting from the next transaction.
<code>pg_backend_get_config_value (<i>pid</i> int4, <i>param_name</i> text, [<i>missing_ok</i> boolean default false], [<i>wait</i> int4 default 1]) → boolean</code>	Returns the current value of the <i>param_name</i> parameter for the backend with the specified PID. If there is no such parameter, the <code>pg_backend_get_config_value</code> function throws an error unless the <i>missing_ok</i> parameter is supplied and is <code>true</code> (in which case NULL is returned). You can specify a timeout in seconds for waiting for the requested parameter value using the optional <i>wait</i> parameter. If the requested parameter value is not returned within the specified timeout, the <code>pg_backend_get_config_value</code> function throws an error.

Function	Description
Example(s)	
	Note that the wait time for the requested parameter value can be indefinite. The target backend will only provide the requested parameter value after completing the current transaction. If the backend is idle, it will not be able to provide the requested parameter value.
<code>pg_backend_load_library (pid int, name text, wait int default 0) → boolean</code>	Loads the library with the provided name for the backend with the specified PID. You can load only one library at a time. When run without superuser privileges, <code>pg_backend_load_library</code> only allows to load libraries located in <code>\$libdir/plugins</code> . If you need to load a library from a different location, run this function on behalf of a superuser.

Functions `pg_backend_set_config` and `pg_backend_load_library` are designed to set configuration parameters and load shared libraries in other sessions, which may be useful for tracing sessions with unexpected behavior. To avoid potential security issues, these functions can only be called by a superuser.

By default, functions `pg_backend_set_config` and `pg_backend_load_library` return `true`, without waiting for response from the other backend. You can use the optional `wait` parameter to specify the time interval, in seconds, for which to wait for the confirmation that the target backend has received the command. If the confirmation is received, the functions return `true`. Otherwise, `false` is returned. Note that the returned value does not reflect the actual result of the operation. If a configuration update fails, an error occurs on the target backend, and the current command on this backend is aborted.

Consider the following examples:

```
SELECT pg_backend_set_config(pg_backend_pid(), 'statement_timeout=10000');
pg_backend_set_config
```

```
-----
t
(1 row)
```

```
SELECT pg_backend_set_config(pg_backend_pid(),
    'statement_timeout=20000
    lock_timeout=10000');
pg_backend_set_config
```

```
-----
t
(1 row)
```

```
SELECT pg_backend_get_config_value(pg_backend_pid(), 'statement_timeout');
pg_backend_get_config_value
```

```
-----
20s
(1 row)
```

```
SELECT pg_backend_get_config_value(pg_backend_pid(), 'lock_timeout');
pg_backend_get_config_value
```

```
-----
10s
(1 row)
```

```
SELECT pg_backend_set_config(pg_backend_pid(), 'fsync=on');
ERROR:  parameter "fsync" cannot be changed now
```

```
SELECT pg_backend_set_config(pg_backend_pid(), 'log_min_messages=INFO');
ERROR:  permission denied to set parameter "log_min_messages"
```

```
SELECT pg_backend_load_library(pg_backend_pid(), 'pgoutput');
pg_backend_load_library
```

```
t
(1 row)
```

9.28.2. Server Signaling Functions

The functions shown in [Table 9.94](#) send control signals to other server processes. Use of these functions is restricted to superusers by default but access may be granted to others using `GRANT`, with noted exceptions.

Each of these functions returns `true` if the signal was successfully sent and `false` if sending the signal failed.

Table 9.94. Server Signaling Functions

Function	Description
<code>pg_cancel_backend (pid integer) → boolean</code>	Cancels the current query of the session whose backend process has the specified process ID. This is also allowed if the calling role is a member of the role whose backend is being canceled or the calling role has privileges of <code>pg_signal_backend</code> , however only superusers can cancel superuser backends.
<code>pg_log_backend_memory_contexts (pid integer) → boolean</code>	Requests to log the memory contexts of the backend with the specified process ID. This function can send the request to backends and auxiliary processes except logger. These memory contexts will be logged at <code>LOG</code> message level. They will appear in the server log based on the log configuration set (see Section 19.8 for more information), but will not be sent to the client regardless of <code>client_min_messages</code> .
<code>pg_reload_conf () → boolean</code>	Causes all processes of the Postgres Pro server to reload their configuration files. (This is initiated by sending a <code>SIGHUP</code> signal to the postmaster process, which in turn sends <code>SIGHUP</code> to each of its children.) You can use the <code>pg_file_settings</code> , <code>pg_hba_file_rules</code> and <code>pg_ident_file_mappings</code> views to check the configuration files for possible errors, before reloading.
<code>pg_rotate_logfile () → boolean</code>	Signals the log-file manager to switch to a new output file immediately. This works only when the built-in log collector is running, since otherwise there is no log-file manager subprocess.
<code>pg_terminate_backend (pid integer, timeout bigint DEFAULT 0) → boolean</code>	Terminates the session whose backend process has the specified process ID. This is also allowed if the calling role is a member of the role whose backend is being terminated or the calling role has privileges of <code>pg_signal_backend</code> , however only superusers can terminate superuser backends. If <code>timeout</code> is not specified or zero, this function returns <code>true</code> whether the process actually terminates or not, indicating only that the sending of the signal was successful. If the <code>timeout</code> is specified (in milliseconds) and greater than zero, the function waits until the process is actually terminated or until the given time has passed. If the process is terminated, the function returns <code>true</code> . On timeout, a warning is emitted and <code>false</code> is returned.

`pg_cancel_backend` and `pg_terminate_backend` send signals (`SIGINT` or `SIGTERM` respectively) to backend processes identified by process ID. The process ID of an active backend can be found from the `pid` column of the `pg_stat_activity` view, or by listing the `postgres` processes on the server (using `ps` on Unix or the Task Manager on Windows). The role of an active backend can be found from the `username` column of the `pg_stat_activity` view.

`pg_log_backend_memory_contexts` can be used to log the memory contexts of a backend process. For example:

```
postgres=# SELECT pg_log_backend_memory_contexts(pg_backend_pid());
pg_log_backend_memory_contexts
-----
t
(1 row)
```

One message for each memory context will be logged. For example:

```
LOG:  logging memory contexts of PID 10377
STATEMENT:  SELECT pg_log_backend_memory_contexts(pg_backend_pid());
LOG:  level: 0; TopMemoryContext: 80800 total in 6 blocks; 14432 free (5 chunks); 66368
      used
LOG:  level: 1; pgstat TabStatusArray lookup hash table: 8192 total in 1 blocks; 1408
      free (0 chunks); 6784 used
LOG:  level: 1; TopTransactionContext: 8192 total in 1 blocks; 7720 free (1 chunks);
      472 used
LOG:  level: 1; RowDescriptionContext: 8192 total in 1 blocks; 6880 free (0 chunks);
      1312 used
LOG:  level: 1; MessageContext: 16384 total in 2 blocks; 5152 free (0 chunks); 11232
      used
LOG:  level: 1; Operator class cache: 8192 total in 1 blocks; 512 free (0 chunks); 7680
      used
LOG:  level: 1; smgr relation table: 16384 total in 2 blocks; 4544 free (3 chunks);
      11840 used
LOG:  level: 1; TransactionAbortContext: 32768 total in 1 blocks; 32504 free (0
      chunks); 264 used
...
LOG:  level: 1; ErrorContext: 8192 total in 1 blocks; 7928 free (3 chunks); 264 used
LOG:  Grand total: 1651920 bytes in 201 blocks; 622360 free (88 chunks); 1029560 used
```

If there are more than 100 child contexts under the same parent, the first 100 child contexts are logged, along with a summary of the remaining contexts. Note that frequent calls to this function could incur significant overhead, because it may generate a large number of log messages.

9.28.3. Backup Control Functions

The functions shown in [Table 9.95](#) assist in making on-line backups. These functions cannot be executed during recovery (except `pg_backup_start`, `pg_backup_stop`, and `pg_wal_lsn_diff`).

For details about proper usage of these functions, see [Section 25.1](#).

Table 9.95. Backup Control Functions

Function	Description
<code>pg_create_restore_point (name text) → pg_lsn</code>	Creates a named marker record in the write-ahead log that can later be used as a recovery target, and returns the corresponding write-ahead log location. The given name can then be used with recovery_target_name to specify the point up to which recovery will proceed. Avoid creating multiple restore points with the same name, since recovery will stop at the first one whose name matches the recovery target. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_current_wal_flush_lsn () → pg_lsn</code>	Returns the current write-ahead log flush location (see notes below).
<code>pg_current_wal_insert_lsn () → pg_lsn</code>	Returns the current write-ahead log insert location (see notes below).
<code>pg_current_wal_lsn () → pg_lsn</code>	Returns the current write-ahead log write location (see notes below).

Function	Description
<code>pg_backup_start (label text [, fast boolean]) → pg_lsn</code>	Prepares the server to begin an on-line backup. The only required parameter is an arbitrary user-defined label for the backup. (Typically this would be the name under which the backup dump file will be stored.) If the optional second parameter is given as <code>true</code> , it specifies executing <code>pg_backup_start</code> as quickly as possible. This forces an immediate checkpoint which will cause a spike in I/O operations, slowing any concurrently executing queries. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_backup_stop ([wait_for_archive boolean]) → record (lsn pg_lsn , label file text , spcmap file text)</code>	Finishes performing an on-line backup. The desired contents of the backup label file and the tablespace map file are returned as part of the result of the function and must be written to files in the backup area. These files must not be written to the live data directory (doing so will cause Postgres Pro to fail to restart in the event of a crash). There is an optional parameter of type <code>boolean</code> . If <code>false</code> , the function will return immediately after the backup is completed, without waiting for WAL to be archived. This behavior is only useful with backup software that independently monitors WAL archiving. Otherwise, WAL required to make the backup consistent might be missing and make the backup useless. By default or when this parameter is <code>true</code> , <code>pg_backup_stop</code> will wait for WAL to be archived when archiving is enabled. (On a standby, this means that it will wait only when <code>archive_mode = always</code> . If write activity on the primary is low, it may be useful to run <code>pg_switch_wal</code> on the primary in order to trigger an immediate segment switch.) When executed on a primary, this function also creates a backup history file in the write-ahead log archive area. The history file includes the label given to <code>pg_backup_start</code> , the starting and ending write-ahead log locations for the backup, and the starting and ending times of the backup. After recording the ending location, the current write-ahead log insertion point is automatically advanced to the next write-ahead log file, so that the ending write-ahead log file can be archived immediately to complete the backup. The result of the function is a single record. The <code>lsn</code> column holds the backup's ending write-ahead log location (which again can be ignored). The second column returns the contents of the backup label file, and the third column returns the contents of the tablespace map file. These must be stored as part of the backup and are required as part of the restore process. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_switch_wal () → pg_lsn</code>	Forces the server to switch to a new write-ahead log file, which allows the current file to be archived (assuming you are using continuous archiving). The result is the ending write-ahead log location plus 1 within the just-completed write-ahead log file. If there has been no write-ahead log activity since the last write-ahead log switch, <code>pg_switch_wal</code> does nothing and returns the start location of the write-ahead log file currently in use. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_walfile_name (lsn pg_lsn) → text</code>	Converts a write-ahead log location to the name of the WAL file holding that location.
<code>pg_walfile_name_offset (lsn pg_lsn) → record (file_name text , file_offset integer)</code>	Converts a write-ahead log location to a WAL file name and byte offset within that file.
<code>pg_split_walfile_name (file_name text) → record (segment_number numeric , timeline_id bigint)</code>	Extracts the sequence number and timeline ID from a WAL file name.
<code>pg_wal_lsn_diff (lsn1 pg_lsn , lsn2 pg_lsn) → numeric</code>	

Function	Description
	Calculates the difference in bytes (<i>lsn1</i> - <i>lsn2</i>) between two write-ahead log locations. This can be used with <code>pg_stat_replication</code> or some of the functions shown in Table 9.95 to get the replication lag.

`pg_current_wal_lsn` displays the current write-ahead log write location in the same format used by the above functions. Similarly, `pg_current_wal_insert_lsn` displays the current write-ahead log insertion location and `pg_current_wal_flush_lsn` displays the current write-ahead log flush location. The insertion location is the “logical” end of the write-ahead log at any instant, while the write location is the end of what has actually been written out from the server's internal buffers, and the flush location is the last location known to be written to durable storage. The write location is the end of what can be examined from outside the server, and is usually what you want if you are interested in archiving partially-complete write-ahead log files. The insertion and flush locations are made available primarily for server debugging purposes. These are all read-only operations and do not require superuser permissions.

You can use `pg_walfile_name_offset` to extract the corresponding write-ahead log file name and byte offset from a `pg_lsn` value. For example:

```
postgres=# SELECT * FROM pg_walfile_name_offset((pg_backup_stop()).lsn);
      file_name      | file_offset
-----+-----
000000010000000000000000D |      4039624
(1 row)
```

Similarly, `pg_walfile_name` extracts just the write-ahead log file name.

`pg_split_walfile_name` is useful to compute a LSN from a file offset and WAL file name, for example:

```
postgres=# \set file_name '000000010000000100C000AB'
postgres=# \set offset 256
postgres=# SELECT '0/0':pg_lsn + pd.segment_number * ps.setting::int + :offset AS lsn
FROM pg_split_walfile_name(:'file_name') pd,
    pg_show_all_settings() ps
WHERE ps.name = 'wal_segment_size';
      lsn
-----
C001/AB000100
(1 row)
```

9.28.4. Recovery Control Functions

The functions shown in [Table 9.96](#) provide information about the current status of a standby server. These functions may be executed both during recovery and in normal running.

Table 9.96. Recovery Information Functions

Function	Description
<code>pg_is_in_recovery</code> () → boolean	Returns true if recovery is still in progress.
<code>pg_last_wal_receive_lsn</code> () → pg_lsn	Returns the last write-ahead log location that has been received and synced to disk by streaming replication. While streaming replication is in progress this will increase monotonically. If recovery has completed then this will remain static at the location of the last WAL record received and synced to disk during recovery. If streaming replication is disabled, or if it has not yet started, the function returns <code>NULL</code> .

Function	Description
<code>pg_last_wal_replay_lsn</code> <code>() → pg_lsn</code>	Returns the last write-ahead log location that has been replayed during recovery. If recovery is still in progress this will increase monotonically. If recovery has completed then this will remain static at the location of the last WAL record applied during recovery. When the server has been started normally without recovery, the function returns <code>NULL</code> .
<code>pg_last_xact_replay_timestamp</code> <code>() → timestamp with time zone</code>	Returns the time stamp of the last transaction replayed during recovery. This is the time at which the commit or abort WAL record for that transaction was generated on the primary. If no transactions have been replayed during recovery, the function returns <code>NULL</code> . Otherwise, if recovery is still in progress this will increase monotonically. If recovery has completed then this will remain static at the time of the last transaction applied during recovery. When the server has been started normally without recovery, the function returns <code>NULL</code> .
<code>pg_get_wal_resource_managers</code> <code>() → setof record (<i>rm_id</i> integer, <i>rm_name</i> text, <i>rm_builtin</i> boolean)</code>	Returns the currently-loaded WAL resource managers in the system. The column <i>rm_builtin</i> indicates whether it's a built-in resource manager, or a custom resource manager loaded by an extension.

The functions shown in [Table 9.97](#) control the progress of recovery. These functions may be executed only during recovery.

Table 9.97. Recovery Control Functions

Function	Description
<code>pg_is_wal_replay_paused</code> <code>() → boolean</code>	Returns true if recovery pause is requested.
<code>pg_get_wal_replay_pause_state</code> <code>() → text</code>	Returns recovery pause state. The return values are <code>not paused</code> if pause is not requested, <code>pause requested</code> if pause is requested but recovery is not yet paused, and <code>paused</code> if the recovery is actually paused.
<code>pg_promote</code> <code>(wait boolean DEFAULT true, wait_seconds integer DEFAULT 60) → boolean</code>	Promotes a standby server to primary status. With <i>wait</i> set to <code>true</code> (the default), the function waits until promotion is completed or <i>wait_seconds</i> seconds have passed, and returns <code>true</code> if promotion is successful and <code>false</code> otherwise. If <i>wait</i> is set to <code>false</code> , the function returns <code>true</code> immediately after sending a <code>SIGUSR1</code> signal to the postmaster to trigger promotion. This function is restricted to superusers by default, but other users can be granted <code>EXECUTE</code> to run the function.
<code>pg_wal_replay_pause</code> <code>() → void</code>	Request to pause recovery. A request doesn't mean that recovery stops right away. If you want a guarantee that recovery is actually paused, you need to check for the recovery pause state returned by <code>pg_get_wal_replay_pause_state()</code> . Note that <code>pg_is_wal_replay_paused()</code> returns whether a request is made. While recovery is paused, no further database changes are applied. If hot standby is active, all new queries will see the same consistent snapshot of the database, and no further query conflicts will be generated until recovery is resumed. This function is restricted to superusers by default, but other users can be granted <code>EXECUTE</code> to run the function.
<code>pg_wal_replay_resume</code> <code>() → void</code>	Restarts recovery if it was paused.

Function	Description
	This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.

`pg_wal_replay_pause` and `pg_wal_replay_resume` cannot be executed while a promotion is ongoing. If a promotion is triggered while recovery is paused, the paused state ends and promotion continues.

If streaming replication is disabled, the paused state may continue indefinitely without a problem. If streaming replication is in progress then WAL records will continue to be received, which will eventually fill available disk space, depending upon the duration of the pause, the rate of WAL generation and available disk space.

9.28.5. Snapshot Synchronization Functions

Postgres Pro allows database sessions to synchronize their snapshots. A *snapshot* determines which data is visible to the transaction that is using the snapshot. Synchronized snapshots are necessary when two or more sessions need to see identical content in the database. If two sessions just start their transactions independently, there is always a possibility that some third transaction commits between the executions of the two `START TRANSACTION` commands, so that one session sees the effects of that transaction and the other does not.

To solve this problem, Postgres Pro allows a transaction to *export* the snapshot it is using. As long as the exporting transaction remains open, other transactions can *import* its snapshot, and thereby be guaranteed that they see exactly the same view of the database that the first transaction sees. But note that any database changes made by any one of these transactions remain invisible to the other transactions, as is usual for changes made by uncommitted transactions. So the transactions are synchronized with respect to pre-existing data, but act normally for changes they make themselves.

Snapshots are exported with the `pg_export_snapshot` function, shown in [Table 9.98](#), and imported with the [SET TRANSACTION](#) command.

Table 9.98. Snapshot Synchronization Functions

Function	Description
<code>pg_export_snapshot () → text</code>	Saves the transaction's current snapshot and returns a <code>text</code> string identifying the snapshot. This string must be passed (outside the database) to clients that want to import the snapshot. The snapshot is available for import only until the end of the transaction that exported it. A transaction can export more than one snapshot, if needed. Note that doing so is only useful in <code>READ COMMITTED</code> transactions, since in <code>REPEATABLE READ</code> and higher isolation levels, transactions use the same snapshot throughout their lifetime. Once a transaction has exported any snapshots, it cannot be prepared with PREPARE TRANSACTION .
<code>pg_log_standby_snapshot () → pg_lsn</code>	Take a snapshot of running transactions and write it to WAL, without having to wait for bg-writer or checkpoint to log one. This is useful for logical decoding on standby, as logical slot creation has to wait until such a record is replayed on the standby.

9.28.6. Replication Management Functions

The functions shown in [Table 9.99](#) are for controlling and interacting with replication features. See [Section 26.2.5](#), [Section 26.2.6](#), and [Chapter 51](#) for information about the underlying features. Use of functions for replication origin is only allowed to the superuser by default, but may be allowed to other users by using the `GRANT` command. Use of functions for replication slots is restricted to superusers and users having `REPLICATION` privilege.

Many of these functions have equivalent commands in the replication protocol; see [Section 56.4](#).

The functions described in [Section 9.28.3](#), [Section 9.28.4](#), and [Section 9.28.5](#) are also relevant for replication.

Table 9.99. Replication Management Functions

Function	Description
<code>pg_create_physical_replication_slot (slot_name name [, immediately_reserve boolean, temporary boolean]) → record (slot_name name, lsn pg_lsn)</code>	Creates a new physical replication slot named <code>slot_name</code> . The optional second parameter, when <code>true</code> , specifies that the LSN for this replication slot be reserved immediately; otherwise the LSN is reserved on first connection from a streaming replication client. Streaming changes from a physical slot is only possible with the streaming-replication protocol — see Section 56.4 . The optional third parameter, <code>temporary</code> , when set to <code>true</code> , specifies that the slot should not be permanently stored to disk and is only meant for use by the current session. Temporary slots are also released upon any error. This function corresponds to the replication protocol command <code>CREATE_REPLICATION_SLOT ... PHYSICAL</code> .
<code>pg_drop_replication_slot (slot_name name) → void</code>	Drops the physical or logical replication slot named <code>slot_name</code> . Same as replication protocol command <code>DROP_REPLICATION_SLOT</code> .
<code>pg_create_logical_replication_slot (slot_name name, plugin name [, temporary boolean, twophase boolean, failover boolean]) → record (slot_name name, lsn pg_lsn)</code>	Creates a new logical (decoding) replication slot named <code>slot_name</code> using the output plugin <code>plugin</code> . The optional third parameter, <code>temporary</code> , when set to <code>true</code> , specifies that the slot should not be permanently stored to disk and is only meant for use by the current session. Temporary slots are also released upon any error. The optional fourth parameter, <code>twophase</code> , when set to <code>true</code> , specifies that the decoding of prepared transactions is enabled for this slot. The optional fifth parameter, <code>failover</code> , when set to <code>true</code> , specifies that this slot is enabled to be synced to the standbys so that logical replication can be resumed after failover. A call to this function has the same effect as the replication protocol command <code>CREATE_REPLICATION_SLOT ... LOGICAL</code> .
<code>pg_copy_physical_replication_slot (src_slot_name name, dst_slot_name name [, temporary boolean]) → record (slot_name name, lsn pg_lsn)</code>	Copies an existing physical replication slot named <code>src_slot_name</code> to a physical replication slot named <code>dst_slot_name</code> . The copied physical slot starts to reserve WAL from the same LSN as the source slot. <code>temporary</code> is optional. If <code>temporary</code> is omitted, the same value as the source slot is used. Copy of an invalidated slot is not allowed.
<code>pg_copy_logical_replication_slot (src_slot_name name, dst_slot_name name [, temporary boolean [, plugin name]]) → record (slot_name name, lsn pg_lsn)</code>	Copies an existing logical replication slot named <code>src_slot_name</code> to a logical replication slot named <code>dst_slot_name</code> , optionally changing the output plugin and persistence. The copied logical slot starts from the same LSN as the source logical slot. Both <code>temporary</code> and <code>plugin</code> are optional; if they are omitted, the values of the source slot are used. The <code>failover</code> option of the source logical slot is not copied and is set to <code>false</code> by default. This is to avoid the risk of being unable to continue logical replication after failover to standby where the slot is being synchronized. Copy of an invalidated slot is not allowed.
<code>pg_logical_slot_get_changes (slot_name name, upto_lsn pg_lsn , upto_nchanges integer, VARIADIC options text[]) → setof record (lsn pg_lsn , xid xid, data text)</code>	Returns changes in the slot <code>slot_name</code> , starting from the point from which changes have been consumed last. If <code>upto_lsn</code> and <code>upto_nchanges</code> are <code>NULL</code> , logical decoding will continue until end of WAL. If <code>upto_lsn</code> is non- <code>NULL</code> , decoding will include only those transactions which commit prior to the specified LSN. If <code>upto_nchanges</code> is non- <code>NULL</code> , decoding will stop when the number of rows produced by decoding exceeds the specified value. Note, however, that the actual number of rows returned may be larger, since this limit is only checked

Function	Description
	after adding the rows produced when decoding each new transaction commit. If the specified slot is a logical failover slot then the function will not return until all physical slots specified in <code>synchronized_standby_slots</code> have confirmed WAL receipt.
<code>pg_logical_slot_peek_changes (slot_name name, upto_lsn pg_lsn , upto_nchanges integer, VARIADIC options text[]) → setof record (lsn pg_lsn , xid xid, data text)</code>	Behaves just like the <code>pg_logical_slot_get_changes()</code> function, except that changes are not consumed; that is, they will be returned again on future calls.
<code>pg_logical_slot_get_binary_changes (slot_name name, upto_lsn pg_lsn , upto_nchanges integer, VARIADIC options text[]) → setof record (lsn pg_lsn , xid xid, data bytea)</code>	Behaves just like the <code>pg_logical_slot_get_changes()</code> function, except that changes are returned as <code>bytea</code> .
<code>pg_logical_slot_peek_binary_changes (slot_name name, upto_lsn pg_lsn , upto_nchanges integer, VARIADIC options text[]) → setof record (lsn pg_lsn , xid xid, data bytea)</code>	Behaves just like the <code>pg_logical_slot_peek_changes()</code> function, except that changes are returned as <code>bytea</code> .
<code>pg_replication_slot_advance (slot_name name, upto_lsn pg_lsn) → record (slot_name name, end_lsn pg_lsn)</code>	Advances the current confirmed position of a replication slot named <code>slot_name</code> . The slot will not be moved backwards, and it will not be moved beyond the current insert location. Returns the name of the slot and the actual position that it was advanced to. The updated slot position information is written out at the next checkpoint if any advancing is done. So in the event of a crash, the slot may return to an earlier position. If the specified slot is a logical failover slot then the function will not return until all physical slots specified in <code>synchronized_standby_slots</code> have confirmed WAL receipt.
<code>pg_replication_origin_create (node_name text) → oid</code>	Creates a replication origin with the given external name, and returns the internal ID assigned to it.
<code>pg_replication_origin_drop (node_name text) → void</code>	Deletes a previously-created replication origin, including any associated replay progress.
<code>pg_replication_origin_oid (node_name text) → oid</code>	Looks up a replication origin by name and returns the internal ID. If no such replication origin is found, <code>NULL</code> is returned.
<code>pg_replication_origin_session_setup (node_name text) → void</code>	Marks the current session as replaying from the given origin, allowing replay progress to be tracked. Can only be used if no origin is currently selected. Use <code>pg_replication_origin_session_reset</code> to undo.
<code>pg_replication_origin_session_reset () → void</code>	Cancels the effects of <code>pg_replication_origin_session_setup()</code> .
<code>pg_replication_origin_session_is_setup () → boolean</code>	Returns true if a replication origin has been selected in the current session.
<code>pg_replication_origin_session_progress (flush boolean) → pg_lsn</code>	Returns the replay location for the replication origin selected in the current session. The parameter <code>flush</code> determines whether the corresponding local transaction will be guaranteed to have been flushed to disk or not.

Function	Description
<code>pg_replication_origin_xact_setup</code>	<code>(origin_lsn pg_lsn , origin_timestamp timestamp with time zone) → void</code> Marks the current transaction as replaying a transaction that has committed at the given LSN and timestamp. Can only be called when a replication origin has been selected using <code>pg_replication_origin_session_setup</code> .
<code>pg_replication_origin_xact_reset</code>	<code>() → void</code> Cancels the effects of <code>pg_replication_origin_xact_setup()</code> .
<code>pg_replication_origin_advance</code>	<code>(node_name text, lsn pg_lsn) → void</code> Sets replication progress for the given node to the given location. This is primarily useful for setting up the initial location, or setting a new location after configuration changes and similar. Be aware that careless use of this function can lead to inconsistently replicated data.
<code>pg_replication_origin_progress</code>	<code>(node_name text, flush boolean) → pg_lsn</code> Returns the replay location for the given replication origin. The parameter <i>flush</i> determines whether the corresponding local transaction will be guaranteed to have been flushed to disk or not.
<code>pg_logical_emit_message</code>	<code>(transactional boolean, prefix text, content text [, flush boolean DEFAULT false]) → pg_lsn</code> <code>pg_logical_emit_message</code> <code>(transactional boolean, prefix text, content bytea [, flush boolean DEFAULT false]) → pg_lsn</code> Emits a logical decoding message. This can be used to pass generic messages to logical decoding plugins through WAL. The <i>transactional</i> parameter specifies if the message should be part of the current transaction, or if it should be written immediately and decoded as soon as the logical decoder reads the record. The <i>prefix</i> parameter is a textual prefix that can be used by logical decoding plugins to easily recognize messages that are interesting for them. The <i>content</i> parameter is the content of the message, given either in text or binary form. The <i>flush</i> parameter (default set to <i>false</i>) controls if the message is immediately flushed to WAL or not. <i>flush</i> has no effect with <i>transactional</i> , as the message's WAL record is flushed along with its transaction.
<code>pg_sync_replication_slots</code>	<code>() → void</code> Synchronize the logical failover replication slots from the primary server to the standby server. This function can only be executed on the standby server. Temporary synced slots, if any, cannot be used for logical decoding and must be dropped after promotion. See Section 50.2.3 for details. Note that this function cannot be executed if <code>sync_replication_slots</code> is enabled and the slotsync worker is already running to perform the synchronization of slots.
<p style="text-align: center;">Caution</p> <p>If, after executing the function, <code>hot_standby_feedback</code> is disabled on the standby or the physical slot configured in <code>primary_slot_name</code> is removed, then it is possible that the necessary rows of the synchronized slot will be removed by the VACUUM process on the primary server, resulting in the synchronized slot becoming invalidated.</p>	

9.28.7. Database Object Management Functions

The functions shown in [Table 9.100](#) calculate the disk space usage of database objects, or assist in presentation or understanding of usage results. `bigint` results are measured in bytes. If an OID that does not represent an existing object is passed to one of these functions, `NULL` is returned.

Table 9.100. Database Object Size Functions

Function	Description
<code>pg_column_size ("any") → integer</code>	Shows the number of bytes used to store any individual data value. If applied directly to a table column value, this reflects any compression that was done.
<code>pg_column_compression ("any") → text</code>	Shows the compression algorithm that was used to compress an individual variable-length value. Returns <code>NULL</code> if the value is not compressed.
<code>pg_column_toast_chunk_id ("any") → oid</code>	Shows the <code>chunk_id</code> of an on-disk TOASTed value. Returns <code>NULL</code> if the value is un-TOASTed or not on-disk. See Section 66.2 for more information about TOAST.
<code>pg_database_size (name) → bigint</code> <code>pg_database_size (oid) → bigint</code>	Computes the total disk space used by the database with the specified name or OID. To use this function, you must have <code>CONNECT</code> privilege on the specified database (which is granted by default) or have privileges of the <code>pg_read_all_stats</code> role.
<code>pg_indexes_size (regclass) → bigint</code>	Computes the total disk space used by indexes attached to the specified table.
<code>pg_relation_size (relation regclass [, fork text]) → bigint</code>	Computes the disk space used by one “fork” of the specified relation. (Note that for most purposes it is more convenient to use the higher-level functions <code>pg_total_relation_size</code> or <code>pg_table_size</code> , which sum the sizes of all forks.) With one argument, this returns the size of the main data fork of the relation. The second argument can be provided to specify which fork to examine: <ul style="list-style-type: none"> • <code>main</code> returns the size of the main data fork of the relation. • <code>fsm</code> returns the size of the Free Space Map (see Section 66.3) associated with the relation. • <code>vm</code> returns the size of the Visibility Map (see Section 66.4) associated with the relation. • <code>init</code> returns the size of the initialization fork, if any, associated with the relation.
<code>pg_temp_relation_size (relation regclass) → bigint</code>	Accepts the OID or name of a temporary table and returns the total size in bytes of <code>main</code> fork of that relation. We extend physical file for temp table when table doesn't fit into temp buffers cache anymore. So total and on-disk size of the temporary table can differ. For other kinds of relations it will not throw any error, just return <code>NULL</code> .
<code>pg_size_bytes (text) → bigint</code>	Converts a size in human-readable format (as returned by <code>pg_size_pretty</code>) into bytes. Valid units are <code>bytes</code> , <code>B</code> , <code>kB</code> , <code>MB</code> , <code>GB</code> , <code>TB</code> , and <code>PB</code> .
<code>pg_size_pretty (bigint) → text</code> <code>pg_size_pretty (numeric) → text</code>	Converts a size in bytes into a more easily human-readable format with size units (bytes, kB, MB, GB, TB, or PB as appropriate). Note that the units are powers of 2 rather than powers of 10, so 1kB is 1024 bytes, 1MB is $1024^2 = 1048576$ bytes, and so on.
<code>pg_table_size (regclass) → bigint</code>	Computes the disk space used by the specified table, excluding indexes (but including its TOAST table if any, free space map, and visibility map).
<code>pg_tablespace_size (name) → bigint</code> <code>pg_tablespace_size (oid) → bigint</code>	

Function	Description
	Computes the total disk space used in the tablespace with the specified name or OID. To use this function, you must have <code>CREATE</code> privilege on the specified tablespace or have privileges of the <code>pg_read_all_stats</code> role, unless it is the default tablespace for the current database.
<code>pg_total_relation_size (regclass) → bigint</code>	Computes the total disk space used by the specified table, including all indexes and TOAST data. The result is equivalent to <code>pg_table_size + pg_indexes_size</code> .

The functions above that operate on tables or indexes accept a `regclass` argument, which is simply the OID of the table or index in the `pg_class` system catalog. You do not have to look up the OID by hand, however, since the `regclass` data type's input converter will do the work for you. See [Section 8.19](#) for details.

The functions shown in [Table 9.101](#) assist in identifying the specific disk files associated with database objects.

Table 9.101. Database Object Location Functions

Function	Description
<code>pg_relation_filenode (relation regclass) → oid</code>	Returns the “filenode” number currently assigned to the specified relation. The filenode is the base component of the file name(s) used for the relation (see Section 66.1 for more information). For most relations the result is the same as <code>pg_class .relfilenode</code> , but for certain system catalogs <code>relfilenode</code> is zero and this function must be used to get the correct value. The function returns <code>NULL</code> if passed a relation that does not have storage, such as a view.
<code>pg_relation_filepath (relation regclass) → text</code>	Returns the entire file path name (relative to the database cluster's data directory, <code>PGDATA</code>) of the relation.
<code>pg_filenode_relation (tablespace oid, filenode oid) → regclass</code>	Returns a relation's OID given the tablespace OID and filenode it is stored under. This is essentially the inverse mapping of <code>pg_relation_filepath</code> . For a relation in the database's default tablespace, the tablespace can be specified as zero. Returns <code>NULL</code> if no relation in the current database is associated with the given values.

[Table 9.102](#) lists functions used to manage collations.

Table 9.102. Collation Management Functions

Function	Description
<code>pg_collation_actual_version (oid) → text</code>	Returns the actual version of the collation object as it is currently installed in the operating system. If this is different from the value in <code>pg_collation .collversion</code> , then objects depending on the collation might need to be rebuilt. See also ALTER COLLATION .
<code>pg_database_collation_actual_version (oid) → text</code>	Returns the actual version of the database's collation as it is currently installed in the operating system. If this is different from the value in <code>pg_database .datcollversion</code> , then objects depending on the collation might need to be rebuilt. See also ALTER DATABASE .
<code>pg_import_system_collations (schema regnamespace) → integer</code>	Adds collations to the system catalog <code>pg_collation</code> based on all the locales it finds in the operating system. This is what <code>initdb</code> uses; see Section 23.2.2 for more details. If additional locales are installed into the operating system later on, this function can be run again to

Function	Description
	add collations for the new locales. Locales that match existing entries in <code>pg_collation</code> will be skipped. (But collation objects based on locales that are no longer present in the operating system are not removed by this function.) The <code>schema</code> parameter would typically be <code>pg_catalog</code> , but that is not a requirement; the collations could be installed into some other schema as well. The function returns the number of new collation objects it created. Use of this function is restricted to superusers.

Table 9.103 lists functions that provide information about the structure of partitioned tables.

Table 9.103. Partitioning Information Functions

Function	Description
<code>pg_partition_tree (regclass) → setof record (relid regclass, parentrelid regclass, isleaf boolean, level integer)</code>	Lists the tables or indexes in the partition tree of the given partitioned table or partitioned index, with one row for each partition. Information provided includes the OID of the partition, the OID of its immediate parent, a boolean value telling if the partition is a leaf, and an integer telling its level in the hierarchy. The level value is 0 for the input table or index, 1 for its immediate child partitions, 2 for their partitions, and so on. Returns no rows if the relation does not exist or is not a partition or partitioned table.
<code>pg_partition_ancestors (regclass) → setof regclass</code>	Lists the ancestor relations of the given partition, including the relation itself. Returns no rows if the relation does not exist or is not a partition or partitioned table.
<code>pg_partition_root (regclass) → regclass</code>	Returns the top-most parent of the partition tree to which the given relation belongs. Returns NULL if the relation does not exist or is not a partition or partitioned table.

For example, to check the total size of the data contained in a partitioned table `measurement`, one could use the following query:

```
SELECT pg_size_pretty (sum(pg_relation_size(relid))) AS total_size
FROM pg_partition_tree ('measurement');
```

9.28.8. Index Maintenance Functions

Table 9.104 shows the functions available for index maintenance tasks. (Note that these maintenance tasks are normally done automatically by autovacuum; use of these functions is only required in special cases.) These functions cannot be executed during recovery. Use of these functions is restricted to superusers and the owner of the given index.

Table 9.104. Index Maintenance Functions

Function	Description
<code>brin_summarize_new_values (index regclass) → integer</code>	Scans the specified BRIN index to find page ranges in the base table that are not currently summarized by the index; for any such range it creates a new summary index tuple by scanning those table pages. Returns the number of new page range summaries that were inserted into the index.
<code>brin_summarize_range (index regclass, blockNumber bigint) → integer</code>	Summarizes the page range covering the given block, if not already summarized. This is like <code>brin_summarize_new_values</code> except that it only processes the page range that covers the given table block number.

Function	Description
<code>brin_desummarize_range (index regclass, blockNumber bigint) → void</code>	Removes the BRIN index tuple that summarizes the page range covering the given table block, if there is one.
<code>gin_clean_pending_list (index regclass) → bigint</code>	Cleans up the “pending” list of the specified GIN index by moving entries in it, in bulk, to the main GIN data structure. Returns the number of pages removed from the pending list. If the argument is a GIN index built with the <code>fastupdate</code> option disabled, no cleanup happens and the result is zero, because the index doesn't have a pending list. See Section 65.4.4.1 and Section 65.4.5 for details about the pending list and <code>fastupdate</code> option.

9.28.9. Generic File Access Functions

The functions shown in [Table 9.105](#) provide native access to files on the machine hosting the server. Only files within the database cluster directory and the `log_directory` can be accessed, unless the user is a superuser or is granted the role `pg_read_server_files`. Use a relative path for files in the cluster directory, and a path matching the `log_directory` configuration setting for log files.

Note that granting users the EXECUTE privilege on `pg_read_file()`, or related functions, allows them the ability to read any file on the server that the database server process can read; these functions bypass all in-database privilege checks. This means that, for example, a user with such access is able to read the contents of the `pg_authid` table where authentication information is stored, as well as read any table data in the database. Therefore, granting access to these functions should be carefully considered.

When granting privilege on these functions, note that the table entries showing optional parameters are mostly implemented as several physical functions with different parameter lists. Privilege must be granted separately on each such function, if it is to be used. `psql`'s `\df` command can be useful to check what the actual function signatures are.

Some of these functions take an optional `missing_ok` parameter, which specifies the behavior when the file or directory does not exist. If `true`, the function returns `NULL` or an empty result set, as appropriate. If `false`, an error is raised. (Failure conditions other than “file not found” are reported as errors in any case.) The default is `false`.

Table 9.105. Generic File Access Functions

Function	Description
<code>pg_ls_dir (dirname text [, missing_ok boolean, include_dot_dirs boolean]) → setof text</code>	Returns the names of all files (and directories and other special files) in the specified directory. The <code>include_dot_dirs</code> parameter indicates whether “.” and “..” are to be included in the result set; the default is to exclude them. Including them can be useful when <code>missing_ok</code> is <code>true</code> , to distinguish an empty directory from a non-existent directory. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_ls_logdir () → setof record (name text, size bigint, modification timestamp with time zone)</code>	Returns the name, size, and last modification time (mtime) of each ordinary file in the server's log directory. Filenames beginning with a dot, directories, and other special files are excluded. This function is restricted to superusers and roles with privileges of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.
<code>pg_ls_waldir () → setof record (name text, size bigint, modification timestamp with time zone)</code>	

Function	Description
	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the server's write-ahead log (WAL) directory. Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and roles with privileges of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_ls_logicalmapdir</code> <code>()</code> \rightarrow <code>setof record (name text, size bigint, modification timestamp with time zone)</code>	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the server's <code>pg_logical/mappings</code> directory. Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and members of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_ls_logicalsnapdir</code> <code>()</code> \rightarrow <code>setof record (name text, size bigint, modification timestamp with time zone)</code>	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the server's <code>pg_logical/snapshots</code> directory. Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and members of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_ls_replslotdir</code> <code>(slot_name text)</code> \rightarrow <code>setof record (name text, size bigint, modification timestamp with time zone)</code>	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the server's <code>pg_replslot/slot_name</code> directory, where <code>slot_name</code> is the name of the replication slot provided as input of the function. Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and members of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_ls_archive_statusdir</code> <code>()</code> \rightarrow <code>setof record (name text, size bigint, modification timestamp with time zone)</code>	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the server's WAL archive status directory (<code>pg_wal/archive_status</code>). Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and members of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_ls_tmpdir</code> <code>([tablespace oid])</code> \rightarrow <code>setof record (name text, size bigint, modification timestamp with time zone)</code>	<p>Returns the name, size, and last modification time (mtime) of each ordinary file in the temporary file directory for the specified <code>tablespace</code>. If <code>tablespace</code> is not provided, the <code>pg_default</code> tablespace is examined. Filenames beginning with a dot, directories, and other special files are excluded.</p> <p>This function is restricted to superusers and members of the <code>pg_monitor</code> role by default, but other users can be granted EXECUTE to run the function.</p>
<code>pg_read_file</code> <code>(filename text [, offset bigint, length bigint] [, missing_ok boolean])</code> \rightarrow <code>text</code>	<p>Returns all or part of a text file, starting at the given byte <code>offset</code>, returning at most <code>length</code> bytes (less if the end of file is reached first). If <code>offset</code> is negative, it is relative to the end of the file. If <code>offset</code> and <code>length</code> are omitted, the entire file is returned. The bytes read from the file are interpreted as a string in the database's encoding; an error is thrown if they are not valid in that encoding.</p> <p>This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.</p>

Function	Description
<code>pg_read_binary_file</code> (<i>filename</i> text [, <i>offset</i> bigint, <i>length</i> bigint] [, <i>missing_ok</i> boolean]) → bytea	Returns all or part of a file. This function is identical to <code>pg_read_file</code> except that it can read arbitrary binary data, returning the result as <code>bytea</code> not <code>text</code> ; accordingly, no encoding checks are performed. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function. In combination with the <code>convert_from</code> function, this function can be used to read a text file in a specified encoding and convert to the database's encoding: SELECT convert_from(pg_read_binary_file('file_in_utf8.txt'), 'UTF8');
<code>pg_stat_file</code> (<i>filename</i> text [, <i>missing_ok</i> boolean]) → record (<i>size</i> bigint, <i>access</i> timestamp with time zone, <i>modification</i> timestamp with time zone, <i>change</i> timestamp with time zone, <i>creation</i> timestamp with time zone, <i>isdir</i> boolean)	Returns a record containing the file's size, last access time stamp, last modification time stamp, last file status change time stamp (Unix platforms only), file creation time stamp (Windows only), and a flag indicating if it is a directory. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.

9.28.10. Advisory Lock Functions

The functions shown in [Table 9.106](#) manage advisory locks. For details about proper use of these functions, see [Section 13.3.5](#).

All these functions are intended to be used to lock application-defined resources, which can be identified either by a single 64-bit key value or two 32-bit key values (note that these two key spaces do not overlap). If another session already holds a conflicting lock on the same resource identifier, the functions will either wait until the resource becomes available, or return a `false` result, as appropriate for the function. Locks can be either shared or exclusive: a shared lock does not conflict with other shared locks on the same resource, only with exclusive locks. Locks can be taken at session level (so that they are held until released or the session ends) or at transaction level (so that they are held until the current transaction ends; there is no provision for manual release). Multiple session-level lock requests stack, so that if the same resource identifier is locked three times there must then be three unlock requests to release the resource in advance of session end.

Table 9.106. Advisory Lock Functions

Function	Description
<code>pg_advisory_lock</code> (<i>key</i> bigint) → void	
<code>pg_advisory_lock</code> (<i>key1</i> integer, <i>key2</i> integer) → void	Obtains an exclusive session-level advisory lock, waiting if necessary.
<code>pg_advisory_lock_shared</code> (<i>key</i> bigint) → void	
<code>pg_advisory_lock_shared</code> (<i>key1</i> integer, <i>key2</i> integer) → void	Obtains a shared session-level advisory lock, waiting if necessary.
<code>pg_advisory_unlock</code> (<i>key</i> bigint) → boolean	
<code>pg_advisory_unlock</code> (<i>key1</i> integer, <i>key2</i> integer) → boolean	Releases a previously-acquired exclusive session-level advisory lock. Returns <code>true</code> if the lock is successfully released. If the lock was not held, <code>false</code> is returned, and in addition, an SQL warning will be reported by the server.
<code>pg_advisory_unlock_all</code> () → void	

Function	Description
	Releases all session-level advisory locks held by the current session. (This function is implicitly invoked at session end, even if the client disconnects ungracefully.)
<code>pg_advisory_unlock_shared (key bigint) → boolean</code> <code>pg_advisory_unlock_shared (key1 integer, key2 integer) → boolean</code>	Releases a previously-acquired shared session-level advisory lock. Returns <code>true</code> if the lock is successfully released. If the lock was not held, <code>false</code> is returned, and in addition, an SQL warning will be reported by the server.
<code>pg_advisory_xact_lock (key bigint) → void</code> <code>pg_advisory_xact_lock (key1 integer, key2 integer) → void</code>	Obtains an exclusive transaction-level advisory lock, waiting if necessary.
<code>pg_advisory_xact_lock_shared (key bigint) → void</code> <code>pg_advisory_xact_lock_shared (key1 integer, key2 integer) → void</code>	Obtains a shared transaction-level advisory lock, waiting if necessary.
<code>pg_try_advisory_lock (key bigint) → boolean</code> <code>pg_try_advisory_lock (key1 integer, key2 integer) → boolean</code>	Obtains an exclusive session-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.
<code>pg_try_advisory_lock_shared (key bigint) → boolean</code> <code>pg_try_advisory_lock_shared (key1 integer, key2 integer) → boolean</code>	Obtains a shared session-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.
<code>pg_try_advisory_xact_lock (key bigint) → boolean</code> <code>pg_try_advisory_xact_lock (key1 integer, key2 integer) → boolean</code>	Obtains an exclusive transaction-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.
<code>pg_try_advisory_xact_lock_shared (key bigint) → boolean</code> <code>pg_try_advisory_xact_lock_shared (key1 integer, key2 integer) → boolean</code>	Obtains a shared transaction-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.

9.28.10.1. Advisory Lock Functions for Distributed System

Postgres Pro Shardman provides a set of functions to get cluster-wide advisory locks. They are similar to regular advisory locks, but are respected on all cluster nodes. For more information, see [Section 70.7](#).

Table 9.107. Advisory Lock Functions for Distributed System

Function	Description
<code>shardman.advisory_xact_lock (key64 BIGINT) → void</code> <code>shardman.advisory_xact_lock (key1 INT, key2 INT) → void</code>	Obtains an exclusive cluster-wide transaction-level advisory lock, waiting if necessary.
<code>shardman.advisory_xact_lock_shared (key64 BIGINT) → void</code> <code>shardman.advisory_xact_lock_shared (key1 INT, key2 INT) → void</code>	

Function	Description
	Obtains a shared cluster-wide transaction-level advisory lock, waiting if necessary.
shardman.try_advisory_xact_lock (key64 BIGINT) → boolean shardman.try_advisory_xact_lock (key1 INT, key2 INT) → boolean	Obtains an exclusive cluster-wide transaction-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.
shardman.try_advisory_xact_lock_shared (key64 BIGINT) → boolean shardman.try_advisory_xact_lock_shared (key1 INT, key2 INT) → boolean	Obtains a shared cluster-wide transaction-level advisory lock if available. This will either obtain the lock immediately and return <code>true</code> , or return <code>false</code> without waiting if the lock cannot be acquired immediately.

9.28.11. Compression Control Functions

The functions shown in [Table 9.108](#) provide information about CFS state and activity and control CFS garbage collection.

Table 9.108. Compression Control Functions

Name	Return Type	Description
cfs_start_gc(n_workers integer)	integer	Starts garbage collection using the specified number of workers. You can only run this function to start the garbage collection manually if background garbage collection is disabled.
cfs_enable_gc(enabled boolean)	boolean	Enables/disables background garbage collection. Alternatively, you can use the cfs_gc configuration variable.
cfs_gc_relation(rel regclass)	integer	Performs garbage collection for a particular table. This function returns the number of processed segments.
cfs_version()	text	Displays the CFS version and the specific compression algorithm used.
cfs_estimate(rel regclass)	float8	Estimates the effect of table compression. Returns the average compression ratio for the first ten blocks of the relation.
cfs_compression_ratio(rel regclass)	float8	Returns the actual compression ratio for all segments of the compressed relation.
cfs_fragmentation(rel regclass)	float8	Returns the average fragmentation ratio of the relation files.
cfs_gc_activity_processed_bytes()	int64	Returns the total size of pages processed by CFS during garbage collection.

Name	Return Type	Description
<code>cfs_gc_activity_processed_pages()</code>	int64	Returns the number of pages processed by CFS during garbage collection.
<code>cfs_gc_activity_processed_files()</code>	int64	Returns the number of files compacted by CFS during garbage collection.
<code>cfs_gc_activity_scanned_files()</code>	int64	Returns the number of files scanned by CFS during garbage collection.

9.28.12. Operation Log Support Functions

An operation log stores information about system events of critical importance, such as an upgrade, execution of [pg_resetwal](#) and so on. This information is not interesting to regular users, but is highly useful for vendor's technical support. Recording to the operation log is only done at the system level (without any actions on the part of the user), and SQL functions are used to read the operation log. The function shown in [Table 9.109](#) reads the operation log.

Table 9.109. Operation Log Support Functions

Function Description
<p><code>pg_operation_log () → setof record (event text, edition text, version text, lsn pg_lsn , last timestampz, count int4)</code></p> <p>Where:</p> <ul style="list-style-type: none"> <i>event</i> — event (operation) type. Can be bootstrap, startup, pg_resetwal , pg_rewind , pg_upgrade or promoted. <i>edition</i> — Postgres Pro edition. Can be vanilla, lc, std, ent or unknown. <i>version</i> — Postgres Pro version. <i>last</i> — date/time of the last occurrence of this event type if events are accumulated or date/time of the event occurrence otherwise. <i>count</i> — number of events of this type if events are accumulated or 1 otherwise. <p>This system function is created automatically for new databases and needs to be created for existing databases, as follows:</p> <pre>CREATE OR REPLACE FUNCTION pg_operation_log(OUT event text, OUT edition text, OUT version text, OUT lsn pg_lsn, OUT last timestampz, OUT count int4) RETURNS SETOF record LANGUAGE INTERNAL STABLE STRICT PARALLEL SAFE AS 'pg_operation_log';</pre>

This is an example of what the `pg_operation_log` function returns:

```
select * from pg_operation_log();
 event | edition | version | lsn | last | count
-----+-----+-----+-----+-----+-----
```

```

startup      | vanilla | 10.22.0 | 0/8000028 | 2022-10-27 23:06:27+03 | 1
pg_upgrade  | 1c      | 15.0.0 | 0/8000028 | 2022-10-27 23:06:27+03 | 1
startup      | 1c      | 15.0.0 | 0/80001F8 | 2022-10-27 23:06:53+03 | 2
(3 rows)

```

Note

If an operation log is empty, before a record for `pg_upgrade`, a record is added for `startup` with the previous *version* having the following specifics: as there is no information in `pg_upgrade` on the patch version number of the previous version for each edition (`1c`, `std`, `ent`), the patch version number (third number in the version number) is set to zero. For an upgrade from `1c` or `vanilla`, it cannot be distinguished whether the edition of the database before upgrade is `vanilla` or `1c`, and `vanilla` edition is used for the `startup` record.

9.28.13. Debugging Functions

The function shown in [Table 9.110](#) can assist you in low-level activities, such as debugging or exploring corrupted Postgres Pro databases.

Table 9.110. Snapshot Synchronization Functions

Name	Return Type	Description
<code>pg_snapshot_any()</code>	<code>void</code>	Sets the current transaction to ignore MVCC rules and see all versions of data.

Use `pg_snapshot_any` with care. Run it in a transaction with isolation level `REPEATABLE READ` or higher, otherwise the specific snapshot will be replaced by a new one by the next query. Only superusers can run this function.

Note

If you created the database cluster using the server version that did not provide this function, execute the command:

```
CREATE FUNCTION pg_snapshot_any() RETURNS void AS 'pg_snapshot_any' LANGUAGE
internal;
```

9.28.14. Troubleshooting Functions

The functions shown in [Table 9.111](#) are used to determine and diagnose issues. For details about proper use of these functions, see [Chapter 34](#).

Table 9.111. Troubleshooting Functions

Function	Description
<code>pg_diagdump (pid integer) → integer</code>	Generates a state dump of the backend process with the specified process ID. This function is typically used for troubleshooting when a process is unresponsive or stuck in an infinite loop. Use of this function is restricted to superusers.

9.29. Trigger Functions

While many uses of triggers involve user-written trigger functions, Postgres Pro provides a few built-in trigger functions that can be used directly in user-defined triggers. These are summarized in [Table 9.112](#).

(Additional built-in trigger functions exist, which implement foreign key constraints and deferred index constraints. Those are not documented here since users need not use them directly.)

For more information about creating triggers, see [CREATE TRIGGER](#).

Table 9.112. Built-In Trigger Functions

Function	Description	Example Usage
<code>suppress_redundant_updates_trigger</code>	<code>() → trigger</code> Suppresses do-nothing update operations. See below for details.	<code>CREATE TRIGGER ... suppress_redundant_updates_trigger()</code>
<code>tsvector_update_trigger</code>	<code>() → trigger</code> Automatically updates a <code>tsvector</code> column from associated plain-text document column(s). The text search configuration to use is specified by name as a trigger argument. See Section 12.4.3 for details.	<code>CREATE TRIGGER ... tsvector_update_trigger(tsvcol, 'pg_catalog.swedish', title, body)</code>
<code>tsvector_update_trigger_column</code>	<code>() → trigger</code> Automatically updates a <code>tsvector</code> column from associated plain-text document column(s). The text search configuration to use is taken from a <code>regconfig</code> column of the table. See Section 12.4.3 for details.	<code>CREATE TRIGGER ... tsvector_update_trigger_column(tsvcol, tsconfigcol, title, body)</code>

The `suppress_redundant_updates_trigger` function, when applied as a row-level `BEFORE UPDATE` trigger, will prevent any update that does not actually change the data in the row from taking place. This overrides the normal behavior which always performs a physical row update regardless of whether or not the data has changed. (This normal behavior makes updates run faster, since no checking is required, and is also useful in certain cases.)

Ideally, you should avoid running updates that don't actually change the data in the record. Redundant updates can cost considerable unnecessary time, especially if there are lots of indexes to alter, and space in dead rows that will eventually have to be vacuumed. However, detecting such situations in client code is not always easy, or even possible, and writing expressions to detect them can be error-prone. An alternative is to use `suppress_redundant_updates_trigger`, which will skip updates that don't change the data. You should use this with care, however. The trigger takes a small but non-trivial time for each record, so if most of the records affected by updates do actually change, use of this trigger will make updates run slower on average.

The `suppress_redundant_updates_trigger` function can be added to a table like this:

```
CREATE TRIGGER z_min_update
BEFORE UPDATE ON tablename
FOR EACH ROW EXECUTE FUNCTION suppress_redundant_updates_trigger();
```

In most cases, you need to fire this trigger last for each row, so that it does not override other triggers that might wish to alter the row. Bearing in mind that triggers fire in name order, you would therefore choose a trigger name that comes after the name of any other trigger you might have on the table. (Hence the “z” prefix in the example.)

9.30. Event Trigger Functions

Postgres Pro provides these helper functions to retrieve information from event triggers.

For more information about event triggers, see [Chapter 41](#).

9.30.1. Capturing Changes at Command End

`pg_event_trigger_ddl_commands () → setof record`

`pg_event_trigger_ddl_commands` returns a list of DDL commands executed by each user action, when invoked in a function attached to a `ddl_command_end` event trigger. If called in any other context, an error is raised. `pg_event_trigger_ddl_commands` returns one row for each base command executed; some commands that are a single SQL sentence may return more than one row. This function returns the following columns:

Name	Type	Description
<code>classid</code>	<code>oid</code>	OID of catalog the object belongs in
<code>objid</code>	<code>oid</code>	OID of the object itself
<code>objsubid</code>	<code>integer</code>	Sub-object ID (e.g., attribute number for a column)
<code>command_tag</code>	<code>text</code>	Command tag
<code>object_type</code>	<code>text</code>	Type of the object
<code>schema_name</code>	<code>text</code>	Name of the schema the object belongs in, if any; otherwise NULL. No quoting is applied.
<code>object_identity</code>	<code>text</code>	Text rendering of the object identity, schema-qualified. Each identifier included in the identity is quoted if necessary.
<code>in_extension</code>	<code>boolean</code>	True if the command is part of an extension script
<code>command</code>	<code>pg_ddl_command</code>	A complete representation of the command, in internal format. This cannot be output directly, but it can be passed to other functions to obtain different pieces of information about the command.

9.30.2. Processing Objects Dropped by a DDL Command

`pg_event_trigger_dropped_objects () → setof record`

`pg_event_trigger_dropped_objects` returns a list of all objects dropped by the command in whose `sql_drop` event it is called. If called in any other context, an error is raised. This function returns the following columns:

Name	Type	Description
<code>classid</code>	<code>oid</code>	OID of catalog the object belonged in
<code>objid</code>	<code>oid</code>	OID of the object itself
<code>objsubid</code>	<code>integer</code>	Sub-object ID (e.g., attribute number for a column)
<code>original</code>	<code>boolean</code>	True if this was one of the root object(s) of the deletion
<code>normal</code>	<code>boolean</code>	True if there was a normal dependency relationship in the de-

Name	Type	Description
		pendency graph leading to this object
is_temporary	boolean	True if this was a temporary object
object_type	text	Type of the object
schema_name	text	Name of the schema the object belonged in, if any; otherwise NULL. No quoting is applied.
object_name	text	Name of the object, if the combination of schema and name can be used as a unique identifier for the object; otherwise NULL. No quoting is applied, and name is never schema-qualified.
object_identity	text	Text rendering of the object identity, schema-qualified. Each identifier included in the identity is quoted if necessary.
address_names	text[]	An array that, together with <code>object_type</code> and <code>address_args</code> , can be used by the <code>pg_get_object_address</code> function to recreate the object address in a remote server containing an identically named object of the same kind.
address_args	text[]	Complement for <code>address_names</code>

The `pg_event_trigger_dropped_objects` function can be used in an event trigger like this:

```
CREATE FUNCTION test_event_trigger_for_drops()
    RETURNS event_trigger LANGUAGE plpgsql AS $$
DECLARE
    obj record;
BEGIN
    FOR obj IN SELECT * FROM pg_event_trigger_dropped_objects()
    LOOP
        RAISE NOTICE '% dropped object: % %.% %',
            tg_tag,
            obj.object_type,
            obj.schema_name,
            obj.object_name,
            obj.object_identity;
    END LOOP;
END;
$$;
CREATE EVENT TRIGGER test_event_trigger_for_drops
    ON sql_drop
    EXECUTE FUNCTION test_event_trigger_for_drops();
```

9.30.3. Handling a Table Rewrite Event

The functions shown in [Table 9.113](#) provide information about a table for which a `table_rewrite` event has just been called. If called in any other context, an error is raised.

Table 9.113. Table Rewrite Information Functions

Function	Description
<code>pg_event_trigger_table_rewrite_oid()</code> \rightarrow oid	Returns the OID of the table about to be rewritten.
<code>pg_event_trigger_table_rewrite_reason()</code> \rightarrow integer	Returns a code explaining the reason(s) for rewriting. The value is a bitmap built from the following values: 1 (the table has changed its persistence), 2 (default value of a column has changed), 4 (a column has a new data type) and 8 (the table access method has changed).

These functions can be used in an event trigger like this:

```
CREATE FUNCTION test_event_trigger_table_rewrite_oid()
RETURNS event_trigger
LANGUAGE plpgsql AS
$$
BEGIN
    RAISE NOTICE 'rewriting table % for reason %',
        pg_event_trigger_table_rewrite_oid()::regclass,
        pg_event_trigger_table_rewrite_reason();
END;
$$;

CREATE EVENT TRIGGER test_table_rewrite_oid
ON table_rewrite
EXECUTE FUNCTION test_event_trigger_table_rewrite_oid();
```

9.31. Statistics Information Functions

Postgres Pro provides a function to inspect complex statistics defined using the `CREATE STATISTICS` command.

9.31.1. Inspecting MCV Lists

`pg_mcv_list_items (pg_mcv_list) \rightarrow setof record`

`pg_mcv_list_items` returns a set of records describing all items stored in a multi-column MCV list. It returns the following columns:

Name	Type	Description
index	integer	index of the item in the MCV list
values	text[]	values stored in the MCV item
nulls	boolean[]	flags identifying NULL values
frequency	double precision	frequency of this MCV item
base_frequency	double precision	base frequency of this MCV item

The `pg_mcv_list_items` function can be used like this:

```
SELECT m.* FROM pg_statistic_ext join pg_statistic_ext_data on (oid = stxoid),
        pg_mcv_list_items(stxdmccv) m WHERE stxname = 'stts';
```

Values of the `pg_mcv_list` type can be obtained only from the `pg_statistic_ext_data.stxdmccv` column.

9.32. Functions for Distributed System

Postgres Pro Shardman provides a set of functions intended for a distributed system. These functions allow you to operate with cluster nodes, manage sharded tables, and broadcast query execution.

`shardman.broadcast_all_sql(statement text)`

Executes *statement* on every replication group.

Warning

The `shardman.broadcast_all_sql` function cannot be executed recursively. An attempt to do so results in an error “Command execution must be initiated by coordinator”.

`shardman.broadcast_query(statement text)`

Functions as `shardman.broadcast_all_sql` and returns an executed SQL *statement* results.

You may optionally set *include_rgid* to `true`, then the resulting tuples will have a number of the node the tuple originated from.

Example with *include_rgid* set to `false`:

```
select shardman.broadcast_query('SELECT relname from pg_class where relkind='f');
broadcast_query
-----
(t_1_fdw)
(t_2_fdw)
(t_0_fdw)
(t_2_fdw)
(t_0_fdw)
(t_1_fdw)
(6 rows)
```

Example with *include_rgid* set to `true`:

```
select shardman.broadcast_query('SELECT relname from pg_class where relkind='f',
include_rgid => true);
broadcast_query
-----
(1,t_1_fdw)
(1,t_2_fdw)
(2,t_0_fdw)
(2,t_2_fdw)
(3,t_0_fdw)
(3,t_1_fdw)
(6 rows)
```

`shardman.broadcast_sql(statement text)`

Executes *statement* on every replication group but the current one.

Warning

The `shardman.broadcast_sql` function cannot be executed recursively. An attempt to do so results in an error “Command execution must be initiated by coordinator”.

`shardman.get_partition_for_value(relid oid, val variadic "any") → shardman.get_partition_for_value_type (rgid int, local_nspname text, local_relname text, remote_nspname text, remote_relname text)`

Finds out which partition of a sharded table with oid *relid* the *val* belongs to. Returns NULL if the sharded table with oid *relid* does not exist. Returns the local schema name and relation names. If the value belongs to a partition stored in another replication group, also returns the remote schema and relation name. Returns only *rgid* if second-level partitioning is used.

Example:

```
select * from shardman.get_partition_for_value('pgbench_branches'::regclass, 20);
rgid | local_nspname |      local_relname      | remote_nspname | remote_relname
-----+-----+-----+-----+-----
+-----+-----+-----+-----+-----
3 | public          | pgbench_branches_17_fdw | public          | pgbench_branches_17
```

```
shardman.global_analyze()
```

Performs cluster-wide analysis of sharded and global tables. First, this function executes `ANALYZE` on all local partitions of sharded tables on each node, then sends this statistics to other nodes. Next, it selects one node per global table and runs `ANALYZE` of this table on the selected node. Gathered statistics is broadcast to all other nodes in the cluster.

Example:

```
select shardman.global_analyze();
```

```
shardman.attach_subpart(relid regclass,  snum int, partition_bound text[])
```

Attaches a previously detached subpartition number *snum* to a locally-partitioned table *relid* as a partition for the values within *partition_bound*. All subpartition tables and foreign tables should already exist. The *partition_bound* parameter is a pair of lower and upper bounds for the partition. If lower and upper bounds are both NULL, the subpartition is attached as the default one.

The operation is performed cluster-wide.

Example:

```
select shardman.attach_subpart('pgbench_history'::regclass, 1,$${'2021-01-01 00:00',
'2022-01-01 00:00'}$$);
```

```
shardman.create_subpart(relid regclass,  snum int, partition_bound text[])
```

Creates a subpartition number *snum* for a locally-partitioned table *relid* as a partition for the values within *partition_bound*. The *partition_bound* parameter is a pair of lower and upper bounds for the partition. If lower and upper bounds are both NULL, the subpartition is created as the default one. If the subpartition number is not specified, it will be selected as the next available partition number.

The operation is performed cluster-wide.

Examples:

```
select shardman.create_subpart('pgbench_history'::regclass, 1, $${'2021-01-01
00:00', '2022-01-01 00:00'}$$);
select shardman.create_subpart('pgbench_history'::regclass, partition_bound:=$
${'2022-01-01 00:00', '2023-01-01 00:00'}$$);
```

```
shardman.detach_subpart(relid regclass,  snum int)
```

Detaches a subpartition number *snum* from a locally-partitioned table *relid*. The partition number can be determined from the `shardman.subparts` view.

The operation is performed cluster-wide.

Example:

```
select shardman.detach_subpart('pgbench_history'::regclass, 1);
```

```
shardman.drop_subpart(relid regclass,  snum int)
```

Drops subpartition number *snum* from locally-partitioned table *relid*. Partition number can be determined from the `shardman.subparts` view.

The operation is performed cluster-wide.

Example:

```
select shardman.drop_subpart('pgbench_history'::regclass, 1);
```

```
shardman.am_coordinator()
```

Returns whether the current session is the query coordinator. This check allows avoiding cases where global and sharded table triggers fire twice, first on the query coordinator, then on the remote nodes when data is modified.

```
SELECT shardman.am_coordinator();
```

```
am_coordinator
```

```
-----
```

```
t
```

```
(1 row)
```

Example of the trigger function checking the query coordinator:

```
CREATE OR REPLACE FUNCTION trg_func() RETURNS TRIGGER
```

```
AS $$
```

```
BEGIN
```

```
IF NOT shardman.am_coordinator() THEN
```

```
    -- exit on non coordinator
```

```
    RETURN NEW;
```

```
END IF;
```

```
-- execute only by coordinator
```

```
RAISE WARNING 'Trigger fired!';
```

```
END
```

```
$$ LANGUAGE plpgsql;
```

Chapter 10. Type Conversion

SQL statements can, intentionally or not, require the mixing of different data types in the same expression. Postgres Pro has extensive facilities for evaluating mixed-type expressions.

In many cases a user does not need to understand the details of the type conversion mechanism. However, implicit conversions done by Postgres Pro can affect the results of a query. When necessary, these results can be tailored by using *explicit* type conversion.

This chapter introduces the Postgres Pro type conversion mechanisms and conventions. Refer to the relevant sections in [Chapter 8](#) and [Chapter 9](#) for more information on specific data types and allowed functions and operators.

10.1. Overview

SQL is a strongly typed language. That is, every data item has an associated data type which determines its behavior and allowed usage. Postgres Pro has an extensible type system that is more general and flexible than other SQL implementations. Hence, most type conversion behavior in Postgres Pro is governed by general rules rather than by ad hoc heuristics. This allows the use of mixed-type expressions even with user-defined types.

The Postgres Pro scanner/parser divides lexical elements into five fundamental categories: integers, non-integer numbers, strings, identifiers, and key words. Constants of most non-numeric types are first classified as strings. The SQL language definition allows specifying type names with strings, and this mechanism can be used in Postgres Pro to start the parser down the correct path. For example, the query:

```
SELECT text 'Origin' AS "label", point '(0,0)' AS "value";
```

```
label | value
-----+-----
Origin | (0,0)
(1 row)
```

has two literal constants, of type `text` and `point`. If a type is not specified for a string literal, then the placeholder type `unknown` is assigned initially, to be resolved in later stages as described below.

There are four fundamental SQL constructs requiring distinct type conversion rules in the Postgres Pro parser:

Function calls

Much of the Postgres Pro type system is built around a rich set of functions. Functions can have one or more arguments. Since Postgres Pro permits function overloading, the function name alone does not uniquely identify the function to be called; the parser must select the right function based on the data types of the supplied arguments.

Operators

Postgres Pro allows expressions with prefix (one-argument) operators, as well as infix (two-argument) operators. Like functions, operators can be overloaded, so the same problem of selecting the right operator exists.

Value Storage

SQL `INSERT` and `UPDATE` statements place the results of expressions into a table. The expressions in the statement must be matched up with, and perhaps converted to, the types of the target columns.

UNION, CASE, and related constructs

Since all query results from a unionized `SELECT` statement must appear in a single set of columns, the types of the results of each `SELECT` clause must be matched up and converted to a uniform set. Similarly, the result expressions of a `CASE` construct must be converted to a common type so that the `CASE` expression as a whole has a known output type. Some other constructs, such as `ARRAY[]`

and the `GREATEST` and `LEAST` functions, likewise require determination of a common type for several subexpressions.

The system catalogs store information about which conversions, or *casts*, exist between which data types, and how to perform those conversions. Additional casts can be added by the user with the `CREATE CAST` command. (This is usually done in conjunction with defining new data types. The set of casts between built-in types has been carefully crafted and is best not altered.)

An additional heuristic provided by the parser allows improved determination of the proper casting behavior among groups of types that have implicit casts. Data types are divided into several basic *type categories*, including `boolean`, `numeric`, `string`, `bitstring`, `datetime`, `timespan`, `geometric`, `network`, and `user-defined`. (For a list see [Table 54.67](#); but note it is also possible to create custom type categories.) Within each category there can be one or more *preferred types*, which are preferred when there is a choice of possible types. With careful selection of preferred types and available implicit casts, it is possible to ensure that ambiguous expressions (those with multiple candidate parsing solutions) can be resolved in a useful way.

All type conversion rules are designed with several principles in mind:

- Implicit conversions should never have surprising or unpredictable outcomes.
- There should be no extra overhead in the parser or executor if a query does not need implicit type conversion. That is, if a query is well-formed and the types already match, then the query should execute without spending extra time in the parser and without introducing unnecessary implicit conversion calls in the query.
- Additionally, if a query usually requires an implicit conversion for a function, and if then the user defines a new function with the correct argument types, the parser should use this new function and no longer do implicit conversion to use the old function.

10.2. Operators

The specific operator that is referenced by an operator expression is determined using the following procedure. Note that this procedure is indirectly affected by the precedence of the operators involved, since that will determine which sub-expressions are taken to be the inputs of which operators. See [Section 4.1.6](#) for more information.

Operator Type Resolution

1. Select the operators to be considered from the `pg_operator` system catalog. If a non-schema-qualified operator name was used (the usual case), the operators considered are those with the matching name and argument count that are visible in the current search path (see [Section 5.10.4](#)). If a qualified operator name was given, only operators in the specified schema are considered.
 - (Optional) If the search path finds multiple operators with identical argument types, only the one appearing earliest in the path is considered. Operators with different argument types are considered on an equal footing regardless of search path position.
2. Check for an operator accepting exactly the input argument types. If one exists (there can be only one exact match in the set of operators considered), use it. Lack of an exact match creates a security hazard when calling, via qualified name ¹ (not typical), any operator found in a schema that permits untrusted users to create objects. In such situations, cast arguments to force an exact match.
 - a. (Optional) If one argument of a binary operator invocation is of the `unknown` type, then assume it is the same type as the other argument for this check. Invocations involving two `unknown` inputs, or a prefix operator with an `unknown` input, will never find a match at this step.
 - b. (Optional) If one argument of a binary operator invocation is of the `unknown` type and the other is of a domain type, next check to see if there is an operator accepting exactly the domain's base type on both sides; if so, use it.

¹ The hazard does not arise with a non-schema-qualified name, because a search path containing schemas that permit untrusted users to create objects is not a [secure schema usage pattern](#).

3. Look for the best match.

- a. Discard candidate operators for which the input types do not match and cannot be converted (using an implicit conversion) to match. `unknown` literals are assumed to be convertible to anything for this purpose. If only one candidate remains, use it; else continue to the next step.
- b. If any input argument is of a domain type, treat it as being of the domain's base type for all subsequent steps. This ensures that domains act like their base types for purposes of ambiguous-operator resolution.
- c. Run through all candidates and keep those with the most exact matches on input types. Keep all candidates if none have exact matches. If only one candidate remains, use it; else continue to the next step.
- d. Run through all candidates and keep those that accept preferred types (of the input data type's type category) at the most positions where type conversion will be required. Keep all candidates if none accept preferred types. If only one candidate remains, use it; else continue to the next step.
- e. If any input arguments are `unknown`, check the type categories accepted at those argument positions by the remaining candidates. At each position, select the `string` category if any candidate accepts that category. (This bias towards string is appropriate since an `unknown`-type literal looks like a string.) Otherwise, if all the remaining candidates accept the same type category, select that category; otherwise fail because the correct choice cannot be deduced without more clues. Now discard candidates that do not accept the selected type category. Furthermore, if any candidate accepts a preferred type in that category, discard candidates that accept non-preferred types for that argument. Keep all candidates if none survive these tests. If only one candidate remains, use it; else continue to the next step.
- f. If there are both `unknown` and `known`-type arguments, and all the `known`-type arguments have the same type, assume that the `unknown` arguments are also of that type, and check which candidates can accept that type at the `unknown`-argument positions. If exactly one candidate passes this test, use it. Otherwise, fail.

Some examples follow.

Example 10.1. Square Root Operator Type Resolution

There is only one square root operator (prefix `|/`) defined in the standard catalog, and it takes an argument of type `double precision`. The scanner assigns an initial type of `integer` to the argument in this query expression:

```
SELECT |/ 40 AS "square root of 40";
square root of 40
-----
6.324555320336759
(1 row)
```

So the parser does a type conversion on the operand and the query is equivalent to:

```
SELECT |/ CAST(40 AS double precision) AS "square root of 40";
```

Example 10.2. String Concatenation Operator Type Resolution

A string-like syntax is used for working with string types and for working with complex extension types. Strings with unspecified type are matched with likely operator candidates.

An example with one unspecified argument:

```
SELECT text 'abc' || 'def' AS "text and unknown";

text and unknown
-----
abcdef
```

```
(1 row)
```

In this case the parser looks to see if there is an operator taking `text` for both arguments. Since there is, it assumes that the second argument should be interpreted as type `text`.

Here is a concatenation of two values of unspecified types:

```
SELECT 'abc' || 'def' AS "unspecified";

 unspecified
-----
 abcdef
(1 row)
```

In this case there is no initial hint for which type to use, since no types are specified in the query. So, the parser looks for all candidate operators and finds that there are candidates accepting both string-category and bit-string-category inputs. Since string category is preferred when available, that category is selected, and then the preferred type for strings, `text`, is used as the specific type to resolve the unknown-type literals as.

Example 10.3. Absolute-Value and Negation Operator Type Resolution

The Postgres Pro operator catalog has several entries for the prefix operator `@`, all of which implement absolute-value operations for various numeric data types. One of these entries is for type `float8`, which is the preferred type in the numeric category. Therefore, Postgres Pro will use that entry when faced with an unknown input:

```
SELECT @ '-4.5' AS "abs";
 abs
-----
 4.5
(1 row)
```

Here the system has implicitly resolved the unknown-type literal as type `float8` before applying the chosen operator. We can verify that `float8` and not some other type was used:

```
SELECT @ '-4.5e500' AS "abs";

ERROR:  "-4.5e500" is out of range for type double precision
```

On the other hand, the prefix operator `~` (bitwise negation) is defined only for integer data types, not for `float8`. So, if we try a similar case with `~`, we get:

```
SELECT ~ '20' AS "negation";

ERROR:  operator is not unique: ~ "unknown"
HINT:   Could not choose a best candidate operator. You might need to add
explicit type casts.
```

This happens because the system cannot decide which of the several possible `~` operators should be preferred. We can help it out with an explicit cast:

```
SELECT ~ CAST('20' AS int8) AS "negation";

 negation
-----
      -21
(1 row)
```

Example 10.4. Array Inclusion Operator Type Resolution

Here is another example of resolving an operator with one known and one unknown input:

```
SELECT array[1,2] <@ '{1,2,3}' as "is subset";
```

```
is subset
-----
t
(1 row)
```

The Postgres Pro operator catalog has several entries for the infix operator `<@`, but the only two that could possibly accept an integer array on the left-hand side are array inclusion (`anyarray <@ anyarray`) and range inclusion (`anyelement <@ anyrange`). Since none of these polymorphic pseudo-types (see [Section 8.21](#)) are considered preferred, the parser cannot resolve the ambiguity on that basis. However, [Step 3.f](#) tells it to assume that the unknown-type literal is of the same type as the other input, that is, integer array. Now only one of the two operators can match, so array inclusion is selected. (Had range inclusion been selected, we would have gotten an error, because the string does not have the right format to be a range literal.)

Example 10.5. Custom Operator on a Domain Type

Users sometimes try to declare operators applying just to a domain type. This is possible but is not nearly as useful as it might seem, because the operator resolution rules are designed to select operators applying to the domain's base type. As an example consider

```
CREATE DOMAIN mytext AS text CHECK(...);
CREATE FUNCTION mytext_eq_text (mytext, text) RETURNS boolean AS ...;
CREATE OPERATOR = (procedure=mytext_eq_text, leftarg=mytext, rightarg=text);
CREATE TABLE mytable (val mytext);

SELECT * FROM mytable WHERE val = 'foo';
```

This query will not use the custom operator. The parser will first see if there is a `mytext = mytext` operator ([Step 2.a](#)), which there is not; then it will consider the domain's base type `text`, and see if there is a `text = text` operator ([Step 2.b](#)), which there is; so it resolves the unknown-type literal as `text` and uses the `text = text` operator. The only way to get the custom operator to be used is to explicitly cast the literal:

```
SELECT * FROM mytable WHERE val = text 'foo';
```

so that the `mytext = text` operator is found immediately according to the exact-match rule. If the best-match rules are reached, they actively discriminate against operators on domain types. If they did not, such an operator would create too many ambiguous-operator failures, because the casting rules always consider a domain as castable to or from its base type, and so the domain operator would be considered usable in all the same cases as a similarly-named operator on the base type.

10.3. Functions

The specific function that is referenced by a function call is determined using the following procedure.

Function Type Resolution

1. Select the functions to be considered from the `pg_proc` system catalog. If a non-schema-qualified function name was used, the functions considered are those with the matching name and argument count that are visible in the current search path (see [Section 5.10.4](#)). If a qualified function name was given, only functions in the specified schema are considered.
 - a. (Optional) If the search path finds multiple functions of identical argument types, only the one appearing earliest in the path is considered. Functions of different argument types are considered on an equal footing regardless of search path position.
 - b. (Optional) If a function is declared with a `VARIADIC` array parameter, and the call does not use the `VARIADIC` keyword, then the function is treated as if the array parameter were replaced by one or more occurrences of its element type, as needed to match the call. After such expansion the function might have effective argument types identical to some non-variadic function. In that case the function appearing earlier in the search path is used, or if the two functions are in the same schema, the non-variadic one is preferred.

This creates a security hazard when calling, via qualified name², a variadic function found in a schema that permits untrusted users to create objects. A malicious user can take control and execute arbitrary SQL functions as though you executed them. Substitute a call bearing the `VARIADIC` keyword, which bypasses this hazard. Calls populating `VARIADIC` "any" parameters often have no equivalent formulation containing the `VARIADIC` keyword. To issue those calls safely, the function's schema must permit only trusted users to create objects.

- c. (Optional) Functions that have default values for parameters are considered to match any call that omits zero or more of the defaultable parameter positions. If more than one such function matches a call, the one appearing earliest in the search path is used. If there are two or more such functions in the same schema with identical parameter types in the non-defaulted positions (which is possible if they have different sets of defaultable parameters), the system will not be able to determine which to prefer, and so an "ambiguous function call" error will result if no better match to the call can be found.

This creates an availability hazard when calling, via qualified name², any function found in a schema that permits untrusted users to create objects. A malicious user can create a function with the name of an existing function, replicating that function's parameters and appending novel parameters having default values. This precludes new calls to the original function. To forestall this hazard, place functions in schemas that permit only trusted users to create objects.

2. Check for a function accepting exactly the input argument types. If one exists (there can be only one exact match in the set of functions considered), use it. Lack of an exact match creates a security hazard when calling, via qualified name², a function found in a schema that permits untrusted users to create objects. In such situations, cast arguments to force an exact match. (Cases involving `unknown` will never find a match at this step.)
3. If no exact match is found, see if the function call appears to be a special type conversion request. This happens if the function call has just one argument and the function name is the same as the (internal) name of some data type. Furthermore, the function argument must be either an `unknown`-type literal, or a type that is binary-coercible to the named data type, or a type that could be converted to the named data type by applying that type's I/O functions (that is, the conversion is either to or from one of the standard string types). When these conditions are met, the function call is treated as a form of `CAST` specification.³
4. Look for the best match.
 - a. Discard candidate functions for which the input types do not match and cannot be converted (using an implicit conversion) to match. `unknown` literals are assumed to be convertible to anything for this purpose. If only one candidate remains, use it; else continue to the next step.
 - b. If any input argument is of a domain type, treat it as being of the domain's base type for all subsequent steps. This ensures that domains act like their base types for purposes of ambiguous-function resolution.
 - c. Run through all candidates and keep those with the most exact matches on input types. Keep all candidates if none have exact matches. If only one candidate remains, use it; else continue to the next step.
 - d. Run through all candidates and keep those that accept preferred types (of the input data type's type category) at the most positions where type conversion will be required. Keep all candidates if none accept preferred types. If only one candidate remains, use it; else continue to the next step.
 - e. If any input arguments are `unknown`, check the type categories accepted at those argument positions by the remaining candidates. At each position, select the `string` category if any candidate accepts that category. (This bias towards string is appropriate since an `unknown`-type literal looks like a string.) Otherwise, if all the remaining candidates accept the same type category,

² The hazard does not arise with a non-schema-qualified name, because a search path containing schemas that permit untrusted users to create objects is not a [secure schema usage pattern](#).

³ The reason for this step is to support function-style cast specifications in cases where there is not an actual cast function. If there is a cast function, it is conventionally named after its output type, and so there is no need to have a special case. See [CREATE CAST](#) for additional commentary.

select that category; otherwise fail because the correct choice cannot be deduced without more clues. Now discard candidates that do not accept the selected type category. Furthermore, if any candidate accepts a preferred type in that category, discard candidates that accept non-preferred types for that argument. Keep all candidates if none survive these tests. If only one candidate remains, use it; else continue to the next step.

- f. If there are both `unknown` and `known-type` arguments, and all the `known-type` arguments have the same type, assume that the `unknown` arguments are also of that type, and check which candidates can accept that type at the `unknown`-argument positions. If exactly one candidate passes this test, use it. Otherwise, fail.

Note that the “best match” rules are identical for operator and function type resolution. Some examples follow.

Example 10.6. Rounding Function Argument Type Resolution

There is only one `round` function that takes two arguments; it takes a first argument of type `numeric` and a second argument of type `integer`. So the following query automatically converts the first argument of type `integer` to `numeric`:

```
SELECT round(4, 4);
```

```
round
-----
4.0000
(1 row)
```

That query is actually transformed by the parser to:

```
SELECT round(CAST (4 AS numeric), 4);
```

Since `numeric` constants with decimal points are initially assigned the type `numeric`, the following query will require no type conversion and therefore might be slightly more efficient:

```
SELECT round(4.0, 4);
```

Example 10.7. Variadic Function Resolution

```
CREATE FUNCTION public.variadic_example(VARIADIC numeric[]) RETURNS int
  LANGUAGE sql AS 'SELECT 1';
CREATE FUNCTION
```

This function accepts, but does not require, the `VARIADIC` keyword. It tolerates both `integer` and `numeric` arguments:

```
SELECT public.variadic_example(0),
       public.variadic_example(0.0),
       public.variadic_example(VARIADIC array[0.0]);
variadic_example | variadic_example | variadic_example
-----+-----+-----
1 | 1 | 1
(1 row)
```

However, the first and second calls will prefer more-specific functions, if available:

```
CREATE FUNCTION public.variadic_example(numeric) RETURNS int
  LANGUAGE sql AS 'SELECT 2';
CREATE FUNCTION

CREATE FUNCTION public.variadic_example(int) RETURNS int
  LANGUAGE sql AS 'SELECT 3';
CREATE FUNCTION

SELECT public.variadic_example(0),
```

```

    public.variadic_example(0.0),
    public.variadic_example(VARIADIC array[0.0]);
variadic_example | variadic_example | variadic_example
-----+-----+-----
          3 |              2 |              1
(1 row)

```

Given the default configuration and only the first function existing, the first and second calls are insecure. Any user could intercept them by creating the second or third function. By matching the argument type exactly and using the `VARIADIC` keyword, the third call is secure.

Example 10.8. Substring Function Type Resolution

There are several `substr` functions, one of which takes types `text` and `integer`. If called with a string constant of unspecified type, the system chooses the candidate function that accepts an argument of the preferred category `string` (namely of type `text`).

```
SELECT substr('1234', 3);
```

```

 substr
-----
      34
(1 row)

```

If the string is declared to be of type `varchar`, as might be the case if it comes from a table, then the parser will try to convert it to become `text`:

```
SELECT substr(varchar '1234', 3);
```

```

 substr
-----
      34
(1 row)

```

This is transformed by the parser to effectively become:

```
SELECT substr(CAST (varchar '1234' AS text), 3);
```

Note

The parser learns from the `pg_cast` catalog that `text` and `varchar` are binary-compatible, meaning that one can be passed to a function that accepts the other without doing any physical conversion. Therefore, no type conversion call is really inserted in this case.

And, if the function is called with an argument of type `integer`, the parser will try to convert that to `text`:

```

SELECT substr(1234, 3);
ERROR:  function substr(integer, integer) does not exist
HINT:  No function matches the given name and argument types. You might need
to add explicit type casts.

```

This does not work because `integer` does not have an implicit cast to `text`. An explicit cast will work, however:

```
SELECT substr(CAST (1234 AS text), 3);
```

```

 substr
-----
      34
(1 row)

```

10.4. Value Storage

Values to be inserted into a table are converted to the destination column's data type according to the following steps.

Value Storage Type Conversion

1. Check for an exact match with the target.
2. Otherwise, try to convert the expression to the target type. This is possible if an *assignment cast* between the two types is registered in the `pg_cast` catalog (see [CREATE CAST](#)). Alternatively, if the expression is an unknown-type literal, the contents of the literal string will be fed to the input conversion routine for the target type.
3. Check to see if there is a sizing cast for the target type. A sizing cast is a cast from that type to itself. If one is found in the `pg_cast` catalog, apply it to the expression before storing into the destination column. The implementation function for such a cast always takes an extra parameter of type `integer`, which receives the destination column's `atttypmod` value (typically its declared length, although the interpretation of `atttypmod` varies for different data types), and it may take a third `boolean` parameter that says whether the cast is explicit or implicit. The cast function is responsible for applying any length-dependent semantics such as size checking or truncation.

Example 10.9. character Storage Type Conversion

For a target column declared as `character(20)` the following statement shows that the stored value is sized correctly:

```
CREATE TABLE vv (v character(20));
INSERT INTO vv SELECT 'abc' || 'def';
SELECT v, octet_length(v) FROM vv;
```

v	octet_length
abcdef	20

(1 row)

What has really happened here is that the two unknown literals are resolved to `text` by default, allowing the `||` operator to be resolved as `text` concatenation. Then the `text` result of the operator is converted to `bpchar` (“blank-padded char”, the internal name of the `character` data type) to match the target column type. (Since the conversion from `text` to `bpchar` is binary-coercible, this conversion does not insert any real function call.) Finally, the sizing function `bpchar(bpchar, integer, boolean)` is found in the system catalog and applied to the operator's result and the stored column length. This type-specific function performs the required length check and addition of padding spaces.

10.5. UNION, CASE, and Related Constructs

SQL `UNION` constructs must match up possibly dissimilar types to become a single result set. The resolution algorithm is applied separately to each output column of a union query. The `INTERSECT` and `EXCEPT` constructs resolve dissimilar types in the same way as `UNION`. Some other constructs, including `CASE`, `ARRAY`, `VALUES`, and the `GREATEST` and `LEAST` functions, use the identical algorithm to match up their component expressions and select a result data type.

Type Resolution for UNION, CASE, and Related Constructs

1. If all inputs are of the same type, and it is not `unknown`, resolve as that type.
2. If any input is of a domain type, treat it as being of the domain's base type for all subsequent steps.⁴
3. If all inputs are of type `unknown`, resolve as type `text` (the preferred type of the string category). Otherwise, `unknown` inputs are ignored for the purposes of the remaining rules.

⁴ Somewhat like the treatment of domain inputs for operators and functions, this behavior allows a domain type to be preserved through a `UNION` or similar construct, so long as the user is careful to ensure that all inputs are implicitly or explicitly of that exact type. Otherwise the domain's base type will be used.

4. If the non-unknown inputs are not all of the same type category, fail.
5. Select the first non-unknown input type as the candidate type, then consider each other non-unknown input type, left to right.⁵ If the candidate type can be implicitly converted to the other type, but not vice-versa, select the other type as the new candidate type. Then continue considering the remaining inputs. If, at any stage of this process, a preferred type is selected, stop considering additional inputs.
6. Convert all inputs to the final candidate type. Fail if there is not an implicit conversion from a given input type to the candidate type.

Some examples follow.

Example 10.10. Type Resolution with Underspecified Types in a Union

```
SELECT text 'a' AS "text" UNION SELECT 'b';
```

```
text
-----
a
b
(2 rows)
```

Here, the unknown-type literal 'b' will be resolved to type `text`.

Example 10.11. Type Resolution in a Simple Union

```
SELECT 1.2 AS "numeric" UNION SELECT 1;
```

```
numeric
-----
      1
     1.2
(2 rows)
```

The literal `1.2` is of type `numeric`, and the integer value `1` can be cast implicitly to `numeric`, so that type is used.

Example 10.12. Type Resolution in a Transposed Union

```
SELECT 1 AS "real" UNION SELECT CAST('2.2' AS REAL);
```

```
real
-----
      1
     2.2
(2 rows)
```

Here, since type `real` cannot be implicitly cast to `integer`, but `integer` can be implicitly cast to `real`, the union result type is resolved as `real`.

Example 10.13. Type Resolution in a Nested Union

```
SELECT NULL UNION SELECT NULL UNION SELECT 1;
```

```
ERROR:  UNION types text and integer cannot be matched
```

This failure occurs because Postgres Pro treats multiple `UNIONS` as a nest of pairwise operations; that is, this input is the same as

```
(SELECT NULL UNION SELECT NULL) UNION SELECT 1;
```

The inner `UNION` is resolved as emitting type `text`, according to the rules given above. Then the outer `UNION` has inputs of types `text` and `integer`, leading to the observed error. The problem can be fixed by ensuring that the leftmost `UNION` has at least one input of the desired result type.

⁵ For historical reasons, `CASE` treats its `ELSE` clause (if any) as the “first” input, with the `THEN` clauses(s) considered after that. In all other cases, “left to right” means the order in which the expressions appear in the query text.

`INTERSECT` and `EXCEPT` operations are likewise resolved pairwise. However, the other constructs described in this section consider all of their inputs in one resolution step.

10.6. SELECT Output Columns

The rules given in the preceding sections will result in assignment of non-unknown data types to all expressions in an SQL query, except for unspecified-type literals that appear as simple output columns of a `SELECT` command. For example, in

```
SELECT 'Hello World';
```

there is nothing to identify what type the string literal should be taken as. In this situation Postgres Pro will fall back to resolving the literal's type as `text`.

When the `SELECT` is one arm of a `UNION` (or `INTERSECT` or `EXCEPT`) construct, or when it appears within `INSERT ... SELECT`, this rule is not applied since rules given in preceding sections take precedence. The type of an unspecified-type literal can be taken from the other `UNION` arm in the first case, or from the destination column in the second case.

`RETURNING` lists are treated the same as `SELECT` output lists for this purpose.

Note

Prior to Postgres Pro 10, this rule did not exist, and unspecified-type literals in a `SELECT` output list were left as type `unknown`. That had assorted bad consequences, so it's been changed.

Chapter 11. Indexes

Indexes are a common way to enhance database performance. An index allows the database server to find and retrieve specific rows much faster than it could do without an index. But indexes also add overhead to the database system as a whole, so they should be used sensibly.

11.1. Introduction

Suppose we have a table similar to this:

```
CREATE TABLE test1 (  
    id integer,  
    content varchar  
);
```

and the application issues many queries of the form:

```
SELECT content FROM test1 WHERE id = constant;
```

With no advance preparation, the system would have to scan the entire `test1` table, row by row, to find all matching entries. If there are many rows in `test1` and only a few rows (perhaps zero or one) that would be returned by such a query, this is clearly an inefficient method. But if the system has been instructed to maintain an index on the `id` column, it can use a more efficient method for locating matching rows. For instance, it might only have to walk a few levels deep into a search tree.

A similar approach is used in most non-fiction books: terms and concepts that are frequently looked up by readers are collected in an alphabetic index at the end of the book. The interested reader can scan the index relatively quickly and flip to the appropriate page(s), rather than having to read the entire book to find the material of interest. Just as it is the task of the author to anticipate the items that readers are likely to look up, it is the task of the database programmer to foresee which indexes will be useful.

The following command can be used to create an index on the `id` column, as discussed:

```
CREATE INDEX test1_id_index ON test1 (id);
```

The name `test1_id_index` can be chosen freely, but you should pick something that enables you to remember later what the index was for.

To remove an index, use the `DROP INDEX` command. Indexes can be added to and removed from tables at any time.

Once an index is created, no further intervention is required: the system will update the index when the table is modified, and it will use the index in queries when it thinks doing so would be more efficient than a sequential table scan. But you might have to run the `ANALYZE` command regularly to update statistics to allow the query planner to make educated decisions. See [Chapter 14](#) for information about how to find out whether an index is used and when and why the planner might choose *not* to use an index.

Indexes can also benefit `UPDATE` and `DELETE` commands with search conditions. Indexes can moreover be used in join searches. Thus, an index defined on a column that is part of a join condition can also significantly speed up queries with joins.

In general, Postgres Pro indexes can be used to optimize queries that contain one or more `WHERE` or `JOIN` clauses of the form

```
indexed-column indexable-operator comparison-value
```

Here, the *indexed-column* is whatever column or expression the index has been defined on. The *indexable-operator* is an operator that is a member of the index's *operator class* for the indexed column. (More details about that appear below.) And the *comparison-value* can be any expression that is not volatile and does not reference the index's table.

In some cases the query planner can extract an indexable clause of this form from another SQL construct. A simple example is that if the original clause was

comparison-value operator indexed-column

then it can be flipped around into indexable form if the original *operator* has a commutator operator that is a member of the index's operator class.

Creating an index on a large table can take a long time. By default, Postgres Pro allows reads (`SELECT` statements) to occur on the table in parallel with index creation, but writes (`INSERT`, `UPDATE`, `DELETE`) are blocked until the index build is finished. In production environments this is often unacceptable. It is possible to allow writes to occur in parallel with index creation, but there are several caveats to be aware of — for more information see [Building Indexes Concurrently](#).

After an index is created, the system has to keep it synchronized with the table. This adds overhead to data manipulation operations. Indexes can also prevent the creation of [heap-only tuples](#). Therefore indexes that are seldom or never used in queries should be removed.

11.2. Index Types

Postgres Pro provides several index types: B-tree, Hash, GiST, SP-GiST, GIN, BRIN. Each index type uses a different algorithm that is best suited to different types of indexable clauses. By default, the `CREATE INDEX` command creates B-tree indexes, which fit the most common situations. The other index types are selected by writing the keyword `USING` followed by the index type name. For example, to create a Hash index:

```
CREATE INDEX name ON table USING HASH (column);
```

11.2.1. B-Tree

B-trees can handle equality and range queries on data that can be sorted into some ordering. In particular, the Postgres Pro query planner will consider using a B-tree index whenever an indexed column is involved in a comparison using one of these operators:

`<` `<=` `=` `>=` `>`

Constructs equivalent to combinations of these operators, such as `BETWEEN` and `IN`, can also be implemented with a B-tree index search. Also, an `IS NULL` or `IS NOT NULL` condition on an index column can be used with a B-tree index.

The optimizer can also use a B-tree index for queries involving the pattern matching operators `LIKE` and `~` if the pattern is a constant and is anchored to the beginning of the string — for example, `col LIKE 'foo%'` or `col ~ '^foo'`, but not `col LIKE '%bar'`. However, if your database does not use the C locale you will need to create the index with a special operator class to support indexing of pattern-matching queries; see [Section 11.10](#) below. It is also possible to use B-tree indexes for `ILIKE` and `~*`, but only if the pattern starts with non-alphabetic characters, i.e., characters that are not affected by upper/lower case conversion.

B-tree indexes can also be used to retrieve data in sorted order. This is not always faster than a simple scan and sort, but it is often helpful.

B-tree indexes are also capable of optimizing “nearest-neighbor” searches, such as

```
SELECT * FROM events ORDER BY event_date <-> date '2017-05-05' LIMIT 10;
```

which finds the ten events closest to a given target date. The ability to do this is again dependent on the particular operator class being used.

11.2.2. Hash

Hash indexes store a 32-bit hash code derived from the value of the indexed column. Hence, such indexes can only handle simple equality comparisons. The query planner will consider using a hash index whenever an indexed column is involved in a comparison using the equal operator:

=

11.2.3. GiST

GiST indexes are not a single kind of index, but rather an infrastructure within which many different indexing strategies can be implemented. Accordingly, the particular operators with which a GiST index can be used vary depending on the indexing strategy (the *operator class*). As an example, the standard distribution of Postgres Pro includes GiST operator classes for several two-dimensional geometric data types, which support indexed queries using these operators:

```
<<    &<    &>    >>    <<|    &<|    |&>    |>>    @>    <@    ~=    &&
```

(See [Section 9.11](#) for the meaning of these operators.) The GiST operator classes included in the standard distribution are documented in [Table 65.1](#). Many other GiST operator classes are available in the `contrib` collection or as separate projects. For more information see [Section 65.2](#).

GiST indexes are also capable of optimizing “nearest-neighbor” searches, such as

```
SELECT * FROM places ORDER BY location <-> point '(101,456)' LIMIT 10;
```

which finds the ten places closest to a given target point. The ability to do this is again dependent on the particular operator class being used. In [Table 65.1](#), operators that can be used in this way are listed in the column “Ordering Operators”.

11.2.4. SP-GiST

SP-GiST indexes, like GiST indexes, offer an infrastructure that supports various kinds of searches. SP-GiST permits implementation of a wide range of different non-balanced disk-based data structures, such as quadrees, k-d trees, and radix trees (tries). As an example, the standard distribution of Postgres Pro includes SP-GiST operator classes for two-dimensional points, which support indexed queries using these operators:

```
<<    >>    ~=    <@    <<|    |>>
```

(See [Section 9.11](#) for the meaning of these operators.) The SP-GiST operator classes included in the standard distribution are documented in [Table 65.2](#). For more information see [Section 65.3](#).

Like GiST, SP-GiST supports “nearest-neighbor” searches. For SP-GiST operator classes that support distance ordering, the corresponding operator is listed in the “Ordering Operators” column in [Table 65.2](#).

11.2.5. GIN

GIN indexes are “inverted indexes” which are appropriate for data values that contain multiple component values, such as arrays. An inverted index contains a separate entry for each component value, and can efficiently handle queries that test for the presence of specific component values.

Like GiST and SP-GiST, GIN can support many different user-defined indexing strategies, and the particular operators with which a GIN index can be used vary depending on the indexing strategy. As an example, the standard distribution of Postgres Pro includes a GIN operator class for arrays, which supports indexed queries using these operators:

```
<@    @>    =    &&
```

(See [Section 9.19](#) for the meaning of these operators.) The GIN operator classes included in the standard distribution are documented in [Table 65.3](#). Many other GIN operator classes are available in the `contrib` collection or as separate projects. For more information see [Section 65.4](#).

11.2.6. BRIN

BRIN indexes (a shorthand for Block Range INdexes) store summaries about the values stored in consecutive physical block ranges of a table. Thus, they are most effective for columns whose values are well-correlated with the physical order of the table rows. Like GiST, SP-GiST and GIN, BRIN can support many different indexing strategies, and the particular operators with which a BRIN index can be used

vary depending on the indexing strategy. For data types that have a linear sort order, the indexed data corresponds to the minimum and maximum values of the values in the column for each block range. This supports indexed queries using these operators:

< <= = >= >

The BRIN operator classes included in the standard distribution are documented in [Table 65.4](#). For more information see [Section 65.5](#).

11.3. Multicolumn Indexes

An index can be defined on more than one column of a table. For example, if you have a table of this form:

```
CREATE TABLE test2 (
    major int,
    minor int,
    name varchar
);
```

(say, you keep your /dev directory in a database...) and you frequently issue queries like:

```
SELECT name FROM test2 WHERE major = constant AND minor = constant;
```

then it might be appropriate to define an index on the columns `major` and `minor` together, e.g.:

```
CREATE INDEX test2_mm_idx ON test2 (major, minor);
```

Currently, only the B-tree, GiST, GIN, and BRIN index types support multiple-key-column indexes. Whether there can be multiple key columns is independent of whether `INCLUDE` columns can be added to the index. Indexes can have up to 32 columns, including `INCLUDE` columns. (This limit can be altered when building Postgres Pro.)

A multicolumn B-tree index can be used with query conditions that involve any subset of the index's columns, but the index is most efficient when there are constraints on the leading (leftmost) columns. The exact rule is that equality constraints on leading columns, plus any inequality constraints on the first column that does not have an equality constraint, will be used to limit the portion of the index that is scanned. Constraints on columns to the right of these columns are checked in the index, so they save visits to the table proper, but they do not reduce the portion of the index that has to be scanned. For example, given an index on (`a`, `b`, `c`) and a query condition `WHERE a = 5 AND b >= 42 AND c < 77`, the index would have to be scanned from the first entry with `a = 5` and `b = 42` up through the last entry with `a = 5`. Index entries with `c >= 77` would be skipped, but they'd still have to be scanned through. This index could in principle be used for queries that have constraints on `b` and/or `c` with no constraint on `a` — but the entire index would have to be scanned, so in most cases the planner would prefer a sequential table scan over using the index.

A multicolumn GiST index can be used with query conditions that involve any subset of the index's columns. Conditions on additional columns restrict the entries returned by the index, but the condition on the first column is the most important one for determining how much of the index needs to be scanned. A GiST index will be relatively ineffective if its first column has only a few distinct values, even if there are many distinct values in additional columns.

A multicolumn GIN index can be used with query conditions that involve any subset of the index's columns. Unlike B-tree or GiST, index search effectiveness is the same regardless of which index column(s) the query conditions use.

A multicolumn BRIN index can be used with query conditions that involve any subset of the index's columns. Like GIN and unlike B-tree or GiST, index search effectiveness is the same regardless of which index column(s) the query conditions use. The only reason to have multiple BRIN indexes instead of one multicolumn BRIN index on a single table is to have a different `pages_per_range` storage parameter.

Of course, each column must be used with operators appropriate to the index type; clauses that involve other operators will not be considered.

Multicolumn indexes should be used sparingly. In most situations, an index on a single column is sufficient and saves space and time. Indexes with more than three columns are unlikely to be helpful unless the usage of the table is extremely stylized. See also [Section 11.5](#) and [Section 11.9](#) for some discussion of the merits of different index configurations.

11.4. Indexes and ORDER BY

In addition to simply finding the rows to be returned by a query, an index may be able to deliver them in a specific sorted order. This allows a query's `ORDER BY` specification to be honored without a separate sorting step. Of the index types currently supported by Postgres Pro, only B-tree can produce sorted output — the other index types return matching rows in an unspecified, implementation-dependent order.

The planner will consider satisfying an `ORDER BY` specification either by scanning an available index that matches the specification, or by scanning the table in physical order and doing an explicit sort. For a query that requires scanning a large fraction of the table, an explicit sort is likely to be faster than using an index because it requires less disk I/O due to following a sequential access pattern. Indexes are more useful when only a few rows need be fetched. An important special case is `ORDER BY` in combination with `LIMIT n`: an explicit sort will have to process all the data to identify the first n rows, but if there is an index matching the `ORDER BY`, the first n rows can be retrieved directly, without scanning the remainder at all.

By default, B-tree indexes store their entries in ascending order with nulls last (table TID is treated as a tiebreaker column among otherwise equal entries). This means that a forward scan of an index on column x produces output satisfying `ORDER BY x` (or more verbosely, `ORDER BY x ASC NULLS LAST`). The index can also be scanned backward, producing output satisfying `ORDER BY x DESC` (or more verbosely, `ORDER BY x DESC NULLS FIRST`, since `NULLS FIRST` is the default for `ORDER BY DESC`).

You can adjust the ordering of a B-tree index by including the options `ASC`, `DESC`, `NULLS FIRST`, and/or `NULLS LAST` when creating the index; for example:

```
CREATE INDEX test2_info_nulls_low ON test2 (info NULLS FIRST);
CREATE INDEX test3_desc_index ON test3 (id DESC NULLS LAST);
```

An index stored in ascending order with nulls first can satisfy either `ORDER BY x ASC NULLS FIRST` or `ORDER BY x DESC NULLS LAST` depending on which direction it is scanned in.

You might wonder why bother providing all four options, when two options together with the possibility of backward scan would cover all the variants of `ORDER BY`. In single-column indexes the options are indeed redundant, but in multicolumn indexes they can be useful. Consider a two-column index on (x, y) : this can satisfy `ORDER BY x, y` if we scan forward, or `ORDER BY x DESC, y DESC` if we scan backward. But it might be that the application frequently needs to use `ORDER BY x ASC, y DESC`. There is no way to get that ordering from a plain index, but it is possible if the index is defined as $(x \text{ ASC}, y \text{ DESC})$ or $(x \text{ DESC}, y \text{ ASC})$.

Obviously, indexes with non-default sort orderings are a fairly specialized feature, but sometimes they can produce tremendous speedups for certain queries. Whether it's worth maintaining such an index depends on how often you use queries that require a special sort ordering.

11.5. Combining Multiple Indexes

A single index scan can only use query clauses that use the index's columns with operators of its operator class and are joined with `AND`. For example, given an index on (a, b) a query condition like `WHERE a = 5 AND b = 6` could use the index, but a query like `WHERE a = 5 OR b = 6` could not directly use the index.

Fortunately, Postgres Pro has the ability to combine multiple indexes (including multiple uses of the same index) to handle cases that cannot be implemented by single index scans. The system can form `AND` and `OR` conditions across several index scans. For example, a query like `WHERE x = 42 OR x = 47 OR x = 53 OR x = 99` could be broken down into four separate scans of an index on x , each scan using one of the query clauses. The results of these scans are then `ORed` together to produce the result. Another example is that if we have separate indexes on x and y , one possible implementation of a query

like `WHERE x = 5 AND y = 6` is to use each index with the appropriate query clause and then AND together the index results to identify the result rows.

To combine multiple indexes, the system scans each needed index and prepares a *bitmap* in memory giving the locations of table rows that are reported as matching that index's conditions. The bitmaps are then ANDed and ORed together as needed by the query. Finally, the actual table rows are visited and returned. The table rows are visited in physical order, because that is how the bitmap is laid out; this means that any ordering of the original indexes is lost, and so a separate sort step will be needed if the query has an `ORDER BY` clause. For this reason, and because each additional index scan adds extra time, the planner will sometimes choose to use a simple index scan even though additional indexes are available that could have been used as well.

In all but the simplest applications, there are various combinations of indexes that might be useful, and the database developer must make trade-offs to decide which indexes to provide. Sometimes multicolumn indexes are best, but sometimes it's better to create separate indexes and rely on the index-combination feature. For example, if your workload includes a mix of queries that sometimes involve only column `x`, sometimes only column `y`, and sometimes both columns, you might choose to create two separate indexes on `x` and `y`, relying on index combination to process the queries that use both columns. You could also create a multicolumn index on `(x, y)`. This index would typically be more efficient than index combination for queries involving both columns, but as discussed in [Section 11.3](#), it would be almost useless for queries involving only `y`, so it should not be the only index. A combination of the multicolumn index and a separate index on `y` would serve reasonably well. For queries involving only `x`, the multicolumn index could be used, though it would be larger and hence slower than an index on `x` alone. The last alternative is to create all three indexes, but this is probably only reasonable if the table is searched much more often than it is updated and all three types of query are common. If one of the types of query is much less common than the others, you'd probably settle for creating just the two indexes that best match the common types.

11.6. Unique Indexes

Indexes can also be used to enforce uniqueness of a column's value, or the uniqueness of the combined values of more than one column.

```
CREATE UNIQUE INDEX name ON table (column [, ...]) [ NULLS [ NOT ] DISTINCT ];
```

Currently, only B-tree indexes can be declared unique.

When an index is declared unique, multiple table rows with equal indexed values are not allowed. By default, null values in a unique column are not considered equal, allowing multiple nulls in the column. The `NULLS NOT DISTINCT` option modifies this and causes the index to treat nulls as equal. A multicolumn unique index will only reject cases where all indexed columns are equal in multiple rows.

Postgres Pro automatically creates a unique index when a unique constraint or primary key is defined for a table. The index covers the columns that make up the primary key or unique constraint (a multicolumn index, if appropriate), and is the mechanism that enforces the constraint.

Note

There's no need to manually create indexes on unique columns; doing so would just duplicate the automatically-created index.

11.7. Indexes on Expressions

An index column need not be just a column of the underlying table, but can be a function or scalar expression computed from one or more columns of the table. This feature is useful to obtain fast access to tables based on the results of computations.

For example, a common way to do case-insensitive comparisons is to use the `lower` function:


```
SELECT * FROM test1 WHERE lower(col1) = 'value';
```

This query can use an index if one has been defined on the result of the `lower(col1)` function:

```
CREATE INDEX test1_lower_col1_idx ON test1 (lower(col1));
```

If we were to declare this index `UNIQUE`, it would prevent creation of rows whose `col1` values differ only in case, as well as rows whose `col1` values are actually identical. Thus, indexes on expressions can be used to enforce constraints that are not definable as simple unique constraints.

As another example, if one often does queries like:

```
SELECT * FROM people WHERE (first_name || ' ' || last_name) = 'John Smith';
```

then it might be worth creating an index like this:

```
CREATE INDEX people_names ON people ((first_name || ' ' || last_name));
```

The syntax of the `CREATE INDEX` command normally requires writing parentheses around index expressions, as shown in the second example. The parentheses can be omitted when the expression is just a function call, as in the first example.

Index expressions are relatively expensive to maintain, because the derived expression(s) must be computed for each row insertion and [non-HOT update](#). However, the index expressions are *not* recomputed during an indexed search, since they are already stored in the index. In both examples above, the system sees the query as just `WHERE indexedcolumn = 'constant'` and so the speed of the search is equivalent to any other simple index query. Thus, indexes on expressions are useful when retrieval speed is more important than insertion and update speed.

11.8. Partial Indexes

A *partial index* is an index built over a subset of a table; the subset is defined by a conditional expression (called the *predicate* of the partial index). The index contains entries only for those table rows that satisfy the predicate. Partial indexes are a specialized feature, but there are several situations in which they are useful.

One major reason for using a partial index is to avoid indexing common values. Since a query searching for a common value (one that accounts for more than a few percent of all the table rows) will not use the index anyway, there is no point in keeping those rows in the index at all. This reduces the size of the index, which will speed up those queries that do use the index. It will also speed up many table update operations because the index does not need to be updated in all cases. [Example 11.1](#) shows a possible application of this idea.

Example 11.1. Setting up a Partial Index to Exclude Common Values

Suppose you are storing web server access logs in a database. Most accesses originate from the IP address range of your organization but some are from elsewhere (say, employees on dial-up connections). If your searches by IP are primarily for outside accesses, you probably do not need to index the IP range that corresponds to your organization's subnet.

Assume a table like this:

```
CREATE TABLE access_log (  
    url varchar,  
    client_ip inet,  
    ...  
);
```

To create a partial index that suits our example, use a command such as this:

```
CREATE INDEX access_log_client_ip_ix ON access_log (client_ip)  
WHERE NOT (client_ip > inet '192.168.100.0' AND  
          client_ip < inet '192.168.100.255');
```

A typical query that can use this index would be:


```
SELECT *
FROM access_log
WHERE url = '/index.html' AND client_ip = inet '212.78.10.32';
```

Here the query's IP address is covered by the partial index. The following query cannot use the partial index, as it uses an IP address that is excluded from the index:

```
SELECT *
FROM access_log
WHERE url = '/index.html' AND client_ip = inet '192.168.100.23';
```

Observe that this kind of partial index requires that the common values be predetermined, so such partial indexes are best used for data distributions that do not change. Such indexes can be recreated occasionally to adjust for new data distributions, but this adds maintenance effort.

Another possible use for a partial index is to exclude values from the index that the typical query workload is not interested in; this is shown in [Example 11.2](#). This results in the same advantages as listed above, but it prevents the “uninteresting” values from being accessed via that index, even if an index scan might be profitable in that case. Obviously, setting up partial indexes for this kind of scenario will require a lot of care and experimentation.

Example 11.2. Setting up a Partial Index to Exclude Uninteresting Values

If you have a table that contains both billed and unbilled orders, where the unbilled orders take up a small fraction of the total table and yet those are the most-accessed rows, you can improve performance by creating an index on just the unbilled rows. The command to create the index would look like this:

```
CREATE INDEX orders_unbilled_index ON orders (order_nr)
WHERE billed is not true;
```

A possible query to use this index would be:

```
SELECT * FROM orders WHERE billed is not true AND order_nr < 10000;
```

However, the index can also be used in queries that do not involve `order_nr` at all, e.g.:

```
SELECT * FROM orders WHERE billed is not true AND amount > 5000.00;
```

This is not as efficient as a partial index on the `amount` column would be, since the system has to scan the entire index. Yet, if there are relatively few unbilled orders, using this partial index just to find the unbilled orders could be a win.

Note that this query cannot use this index:

```
SELECT * FROM orders WHERE order_nr = 3501;
```

The order 3501 might be among the billed or unbilled orders.

[Example 11.2](#) also illustrates that the indexed column and the column used in the predicate do not need to match. Postgres Pro supports partial indexes with arbitrary predicates, so long as only columns of the table being indexed are involved. However, keep in mind that the predicate must match the conditions used in the queries that are supposed to benefit from the index. To be precise, a partial index can be used in a query only if the system can recognize that the `WHERE` condition of the query mathematically implies the predicate of the index. Postgres Pro does not have a sophisticated theorem prover that can recognize mathematically equivalent expressions that are written in different forms. (Not only is such a general theorem prover extremely difficult to create, it would probably be too slow to be of any real use.) The system can recognize simple inequality implications, for example “ $x < 1$ ” implies “ $x < 2$ ”; otherwise the predicate condition must exactly match part of the query's `WHERE` condition or the index will not be recognized as usable. Matching takes place at query planning time, not at run time. As a result, parameterized query clauses do not work with a partial index. For example a prepared query with a parameter might specify “ $x < ?$ ” which will never imply “ $x < 2$ ” for all possible values of the parameter.

A third possible use for partial indexes does not require the index to be used in queries at all. The idea here is to create a unique index over a subset of a table, as in [Example 11.3](#). This enforces uniqueness among the rows that satisfy the index predicate, without constraining those that do not.

Example 11.3. Setting up a Partial Unique Index

Suppose that we have a table describing test outcomes. We wish to ensure that there is only one “successful” entry for a given subject and target combination, but there might be any number of “unsuccessful” entries. Here is one way to do it:

```
CREATE TABLE tests (
    subject text,
    target text,
    success boolean,
    ...
);

CREATE UNIQUE INDEX tests_success_constraint ON tests (subject, target)
    WHERE success;
```

This is a particularly efficient approach when there are few successful tests and many unsuccessful ones. It is also possible to allow only one null in a column by creating a unique partial index with an `IS NULL` restriction.

Finally, a partial index can also be used to override the system's query plan choices. Also, data sets with peculiar distributions might cause the system to use an index when it really should not. In that case the index can be set up so that it is not available for the offending query. Normally, Postgres Pro makes reasonable choices about index usage (e.g., it avoids them when retrieving common values, so the earlier example really only saves index size, it is not required to avoid index usage), and grossly incorrect plan choices are cause for a bug report.

Keep in mind that setting up a partial index indicates that you know at least as much as the query planner knows, in particular you know when an index might be profitable. Forming this knowledge requires experience and understanding of how indexes in Postgres Pro work. In most cases, the advantage of a partial index over a regular index will be minimal. There are cases where they are quite counterproductive, as in [Example 11.4](#).

Example 11.4. Do Not Use Partial Indexes as a Substitute for Partitioning

You might be tempted to create a large set of non-overlapping partial indexes, for example

```
CREATE INDEX mytable_cat_1 ON mytable (data) WHERE category = 1;
CREATE INDEX mytable_cat_2 ON mytable (data) WHERE category = 2;
CREATE INDEX mytable_cat_3 ON mytable (data) WHERE category = 3;
...
CREATE INDEX mytable_cat_N ON mytable (data) WHERE category = N;
```

This is a bad idea! Almost certainly, you'll be better off with a single non-partial index, declared like

```
CREATE INDEX mytable_cat_data ON mytable (category, data);
```

(Put the category column first, for the reasons described in [Section 11.3](#).) While a search in this larger index might have to descend through a couple more tree levels than a search in a smaller index, that's almost certainly going to be cheaper than the planner effort needed to select the appropriate one of the partial indexes. The core of the problem is that the system does not understand the relationship among the partial indexes, and will laboriously test each one to see if it's applicable to the current query.

If your table is large enough that a single index really is a bad idea, you should look into using partitioning instead (see [Section 5.12](#)). With that mechanism, the system does understand that the tables and indexes are non-overlapping, so far better performance is possible.

More information about partial indexes can be found in [ston89b](#), [olson93](#), and [seshadri95](#).

11.9. Index-Only Scans and Covering Indexes

All indexes in Postgres Pro are *secondary* indexes, meaning that each index is stored separately from the table's main data area (which is called the table's *heap* in Postgres Pro terminology). This means

that in an ordinary index scan, each row retrieval requires fetching data from both the index and the heap. Furthermore, while the index entries that match a given indexable `WHERE` condition are usually close together in the index, the table rows they reference might be anywhere in the heap. The heap-access portion of an index scan thus involves a lot of random access into the heap, which can be slow, particularly on traditional rotating media. (As described in [Section 11.5](#), bitmap scans try to alleviate this cost by doing the heap accesses in sorted order, but that only goes so far.)

To solve this performance problem, Postgres Pro supports *index-only scans*, which can answer queries from an index alone without any heap access. The basic idea is to return values directly out of each index entry instead of consulting the associated heap entry. There are two fundamental restrictions on when this method can be used:

1. The index type must support index-only scans. B-tree indexes always do. GiST and SP-GiST indexes support index-only scans for some operator classes but not others. Other index types have no support. The underlying requirement is that the index must physically store, or else be able to reconstruct, the original data value for each index entry. As a counterexample, GIN indexes cannot support index-only scans because each index entry typically holds only part of the original data value.
2. The query must reference only columns stored in the index. For example, given an index on columns `x` and `y` of a table that also has a column `z`, these queries could use index-only scans:

```
SELECT x, y FROM tab WHERE x = 'key';
SELECT x FROM tab WHERE x = 'key' AND y < 42;
```

but these queries could not:

```
SELECT x, z FROM tab WHERE x = 'key';
SELECT x FROM tab WHERE x = 'key' AND z < 42;
```

(Expression indexes and partial indexes complicate this rule, as discussed below.)

If these two fundamental requirements are met, then all the data values required by the query are available from the index, so an index-only scan is physically possible. But there is an additional requirement for any table scan in Postgres Pro: it must verify that each retrieved row be “visible” to the query’s MVCC snapshot, as discussed in [Chapter 13](#). Visibility information is not stored in index entries, only in heap entries; so at first glance it would seem that every row retrieval would require a heap access anyway. And this is indeed the case, if the table row has been modified recently. However, for seldom-changing data there is a way around this problem. Postgres Pro tracks, for each page in a table’s heap, whether all rows stored in that page are old enough to be visible to all current and future transactions. This information is stored in a bit in the table’s *visibility map*. An index-only scan, after finding a candidate index entry, checks the visibility map bit for the corresponding heap page. If it’s set, the row is known visible and so the data can be returned with no further work. If it’s not set, the heap entry must be visited to find out whether it’s visible, so no performance advantage is gained over a standard index scan. Even in the successful case, this approach trades visibility map accesses for heap accesses; but since the visibility map is four orders of magnitude smaller than the heap it describes, far less physical I/O is needed to access it. In most situations the visibility map remains cached in memory all the time.

In short, while an index-only scan is possible given the two fundamental requirements, it will be a win only if a significant fraction of the table’s heap pages have their all-visible map bits set. But tables in which a large fraction of the rows are unchanging are common enough to make this type of scan very useful in practice.

To make effective use of the index-only scan feature, you might choose to create a *covering index*, which is an index specifically designed to include the columns needed by a particular type of query that you run frequently. Since queries typically need to retrieve more columns than just the ones they search on, Postgres Pro allows you to create an index in which some columns are just “payload” and are not part of the search key. This is done by adding an `INCLUDE` clause listing the extra columns. For example, if you commonly run queries like

```
SELECT y FROM tab WHERE x = 'key';
```

the traditional approach to speeding up such queries would be to create an index on `x` only. However, an index defined as

```
CREATE INDEX tab_x_y ON tab(x) INCLUDE (y);
```

could handle these queries as index-only scans, because `y` can be obtained from the index without visiting the heap.

Because column `y` is not part of the index's search key, it does not have to be of a data type that the index can handle; it's merely stored in the index and is not interpreted by the index machinery. Also, if the index is a unique index, that is

```
CREATE UNIQUE INDEX tab_x_y ON tab(x) INCLUDE (y);
```

the uniqueness condition applies to just column `x`, not to the combination of `x` and `y`. (An `INCLUDE` clause can also be written in `UNIQUE` and `PRIMARY KEY` constraints, providing alternative syntax for setting up an index like this.)

It's wise to be conservative about adding non-key payload columns to an index, especially wide columns. If an index tuple exceeds the maximum size allowed for the index type, data insertion will fail. In any case, non-key columns duplicate data from the index's table and bloat the size of the index, thus potentially slowing searches. And remember that there is little point in including payload columns in an index unless the table changes slowly enough that an index-only scan is likely to not need to access the heap. If the heap tuple must be visited anyway, it costs nothing more to get the column's value from there. Other restrictions are that expressions are not currently supported as included columns, and that only B-tree, GiST and SP-GiST indexes currently support included columns.

Before Postgres Pro had the `INCLUDE` feature, people sometimes made covering indexes by writing the payload columns as ordinary index columns, that is writing

```
CREATE INDEX tab_x_y ON tab(x, y);
```

even though they had no intention of ever using `y` as part of a `WHERE` clause. This works fine as long as the extra columns are trailing columns; making them be leading columns is unwise for the reasons explained in [Section 11.3](#). However, this method doesn't support the case where you want the index to enforce uniqueness on the key column(s).

Suffix truncation always removes non-key columns from upper B-Tree levels. As payload columns, they are never used to guide index scans. The truncation process also removes one or more trailing key column(s) when the remaining prefix of key column(s) happens to be sufficient to describe tuples on the lowest B-Tree level. In practice, covering indexes without an `INCLUDE` clause often avoid storing columns that are effectively payload in the upper levels. However, explicitly defining payload columns as non-key columns *reliably* keeps the tuples in upper levels small.

In principle, index-only scans can be used with expression indexes. For example, given an index on `f(x)` where `x` is a table column, it should be possible to execute

```
SELECT f(x) FROM tab WHERE f(x) < 1;
```

as an index-only scan; and this is very attractive if `f()` is an expensive-to-compute function. However, Postgres Pro's planner is currently not very smart about such cases. It considers a query to be potentially executable by index-only scan only when all *columns* needed by the query are available from the index. In this example, `x` is not needed except in the context `f(x)`, but the planner does not notice that and concludes that an index-only scan is not possible. If an index-only scan seems sufficiently worthwhile, this can be worked around by adding `x` as an included column, for example

```
CREATE INDEX tab_f_x ON tab (f(x)) INCLUDE (x);
```

An additional caveat, if the goal is to avoid recalculating `f(x)`, is that the planner won't necessarily match uses of `f(x)` that aren't in indexable `WHERE` clauses to the index column. It will usually get this right in simple queries such as shown above, but not in queries that involve joins. These deficiencies may be remedied in future versions of Postgres Pro.

Partial indexes also have interesting interactions with index-only scans. Consider the partial index shown in [Example 11.3](#):

```
CREATE UNIQUE INDEX tests_success_constraint ON tests (subject, target)
WHERE success;
```

In principle, we could do an index-only scan on this index to satisfy a query like

```
SELECT target FROM tests WHERE subject = 'some-subject' AND success;
```

But there's a problem: the `WHERE` clause refers to `success` which is not available as a result column of the index. Nonetheless, an index-only scan is possible because the plan does not need to recheck that part of the `WHERE` clause at run time: all entries found in the index necessarily have `success = true` so this need not be explicitly checked in the plan. Postgres Pro versions 9.6 and later will recognize such cases and allow index-only scans to be generated, but older versions will not.

11.10. Operator Classes and Operator Families

An index definition can specify an *operator class* for each column of an index.

```
CREATE INDEX name ON table (column opclass [ ( opclass_options ) ] [sort options]
[, ...]);
```

The operator class identifies the operators to be used by the index for that column. For example, a B-tree index on the type `int4` would use the `int4_ops` class; this operator class includes comparison functions for values of type `int4`. In practice the default operator class for the column's data type is usually sufficient. The main reason for having operator classes is that for some data types, there could be more than one meaningful index behavior. For example, we might want to sort a complex-number data type either by absolute value or by real part. We could do this by defining two operator classes for the data type and then selecting the proper class when making an index. The operator class determines the basic sort ordering (which can then be modified by adding sort options `COLLATE`, `ASC/DESC` and/or `NULLS FIRST/NULLS LAST`).

There are also some built-in operator classes besides the default ones:

- The operator classes `text_pattern_ops`, `varchar_pattern_ops`, and `bpchar_pattern_ops` support B-tree indexes on the types `text`, `varchar`, and `char` respectively. The difference from the default operator classes is that the values are compared strictly character by character rather than according to the locale-specific collation rules. This makes these operator classes suitable for use by queries involving pattern matching expressions (`LIKE` or POSIX regular expressions) when the database does not use the standard “C” locale. As an example, you might index a `varchar` column like this:

```
CREATE INDEX test_index ON test_table (col varchar_pattern_ops);
```

Note that you should also create an index with the default operator class if you want queries involving ordinary `<`, `<=`, `>`, or `>=` comparisons to use an index. Such queries cannot use the `xxx_pattern_ops` operator classes. (Ordinary equality comparisons can use these operator classes, however.) It is possible to create multiple indexes on the same column with different operator classes. If you do use the C locale, you do not need the `xxx_pattern_ops` operator classes, because an index with the default operator class is usable for pattern-matching queries in the C locale.

The following query shows all defined operator classes:

```
SELECT am.amname AS index_method,
       opc.opcname AS opclass_name,
       opc.opcintype::regtype AS indexed_type,
       opc.opcdefault AS is_default
FROM pg_am am, pg_opclass opc
WHERE opc.opcmethod = am.oid
ORDER BY index_method, opclass_name;
```

An operator class is actually just a subset of a larger structure called an *operator family*. In cases where several data types have similar behaviors, it is frequently useful to define cross-data-type operators and allow these to work with indexes. To do this, the operator classes for each of the types must be grouped

into the same operator family. The cross-type operators are members of the family, but are not associated with any single class within the family.

This expanded version of the previous query shows the operator family each operator class belongs to:

```
SELECT am.amname AS index_method,
       opc.opcname AS opclass_name,
       opf.opfname AS opfamily_name,
       opc.opcintype::regtype AS indexed_type,
       opc.opcdefault AS is_default
FROM pg_am am, pg_opclass opc, pg_opfamily opf
WHERE opc.opcmethod = am.oid AND
       opc.opcfamily = opf.oid
ORDER BY index_method, opclass_name;
```

This query shows all defined operator families and all the operators included in each family:

```
SELECT am.amname AS index_method,
       opf.opfname AS opfamily_name,
       amop.amopopr::regoperator AS opfamily_operator
FROM pg_am am, pg_opfamily opf, pg_amop amop
WHERE opf.opfmethod = am.oid AND
       amop.amopfamilly = opf.oid
ORDER BY index_method, opfamily_name, opfamily_operator;
```

Tip

psql has commands `\dAc`, `\dAf`, and `\dAo`, which provide slightly more sophisticated versions of these queries.

11.11. Indexes and Collations

An index can support only one collation per index column. If multiple collations are of interest, multiple indexes may be needed.

Consider these statements:

```
CREATE TABLE test1c (
    id integer,
    content varchar COLLATE "x"
);
```

```
CREATE INDEX test1c_content_index ON test1c (content);
```

The index automatically uses the collation of the underlying column. So a query of the form

```
SELECT * FROM test1c WHERE content > constant;
```

could use the index, because the comparison will by default use the collation of the column. However, this index cannot accelerate queries that involve some other collation. So if queries of the form, say,

```
SELECT * FROM test1c WHERE content > constant COLLATE "y";
```

are also of interest, an additional index could be created that supports the "y" collation, like this:

```
CREATE INDEX test1c_content_y_index ON test1c (content COLLATE "y");
```

11.12. Examining Index Usage

Although indexes in Postgres Pro do not need maintenance or tuning, it is still important to check which indexes are actually used by the real-life query workload. Examining index usage for an individual query

is done with the [EXPLAIN](#) command; its application for this purpose is illustrated in [Section 14.1](#). It is also possible to gather overall statistics about index usage in a running server, as described in [Section 27.2](#).

It is difficult to formulate a general procedure for determining which indexes to create. There are a number of typical cases that have been shown in the examples throughout the previous sections. A good deal of experimentation is often necessary. The rest of this section gives some tips for that:

- Always run [ANALYZE](#) first. This command collects statistics about the distribution of the values in the table. This information is required to estimate the number of rows returned by a query, which is needed by the planner to assign realistic costs to each possible query plan. In absence of any real statistics, some default values are assumed, which are almost certain to be inaccurate. Examining an application's index usage without having run [ANALYZE](#) is therefore a lost cause. See [Section 24.1.3](#) and [Section 24.1.6](#) for more information.
- Use real data for experimentation. Using test data for setting up indexes will tell you what indexes you need for the test data, but that is all.

It is especially fatal to use very small test data sets. While selecting 1000 out of 100000 rows could be a candidate for an index, selecting 1 out of 100 rows will hardly be, because the 100 rows probably fit within a single disk page, and there is no plan that can beat sequentially fetching 1 disk page.

Also be careful when making up test data, which is often unavoidable when the application is not yet in production. Values that are very similar, completely random, or inserted in sorted order will skew the statistics away from the distribution that real data would have.

- When indexes are not used, it can be useful for testing to force their use. There are run-time parameters that can turn off various plan types (see [Section 19.7.1](#)). For instance, turning off sequential scans (`enable_seqscan`) and nested-loop joins (`enable_nestloop`), which are the most basic plans, will force the system to use a different plan. If the system still chooses a sequential scan or nested-loop join then there is probably a more fundamental reason why the index is not being used; for example, the query condition does not match the index. (What kind of query can use what kind of index is explained in the previous sections.)
- If forcing index usage does use the index, then there are two possibilities: Either the system is right and using the index is indeed not appropriate, or the cost estimates of the query plans are not reflecting reality. So you should time your query with and without indexes. The `EXPLAIN ANALYZE` command can be useful here.
- If it turns out that the cost estimates are wrong, there are, again, two possibilities. The total cost is computed from the per-row costs of each plan node times the selectivity estimate of the plan node. The costs estimated for the plan nodes can be adjusted via run-time parameters (described in [Section 19.7.2](#)). An inaccurate selectivity estimate is due to insufficient statistics. It might be possible to improve this by tuning the statistics-gathering parameters (see [ALTER TABLE](#)).

If you do not succeed in adjusting the costs to be more appropriate, then you might have to resort to forcing index usage explicitly. You might also want to contact the Postgres Pro developers to examine the issue.

11.13. Using k-NN Algorithm for Optimizing Queries

Postgres Pro Shardman can optimize k-nearest neighbors (k-NN) search for B-tree, GiST, and SP-GiST index types. With these index types, you can use ordering operators that return the rows in the order defined by the `ORDER BY` clause of the following type:

```
ORDER BY index_column operator constant
```

where *operator* is the supported ordering operator, such as `<->`. In this case, Postgres Pro Shardman uses a special scan strategy for k-NN search.

For B-tree indexes, you can use only one ordering operator on the first index column. GiST and SP-GiST do not limit the number of ordering operators you can use for k-NN search. For multicolumn GiST

indexes, you can also use multiple ordering operators for each column included into the multicolumn index. For example:

```
SELECT * FROM points ORDER BY column1 <-> '(0, 0)', column2 <-> '(1, 1)'
```

where *column1* and *column2* are index columns of the multicolumn GiST index.

You can run k-NN search for all built-in and custom data types that support the concept of distance. In database queries, k-NN search can provide performance benefits for various use cases: geographic information systems (GIS), classification problems (majority voting), clusterization, looking for similar objects or nearest events, detecting duplicates and typos, trigram matching, multimedia search, and more.

To natively implement k-NN search, Postgres Pro Shardman uses distance-based priority queue for index traversal. Depending on the index type, the search algorithm will differ:

- For GiST and SP-GiST index types, the priority queue is based on the pairing heap that can be effectively balanced. Inner nodes of the tree represent the minimum bounding rectangle (MBR) for their children. Leaf nodes store actual tuple entries. Postgres Pro Shardman uses a heuristic to maintain the order by distance between the queue entries and the target point defined by the search argument. For inner nodes, the minimal distance for the child node is calculated. For leaf nodes, the distance to the tuple entry is used. In the tree nodes, leaf nodes are stored before the inner nodes in the list, so tuple entries are processed first. Each time a result is returned from the queue, it is guaranteed to be the closest one to the target point among the remaining entries. Thus, there is no need to scan all nodes. Once the required k entries are extracted from the queue, search can stop.
- B-tree indexes use a simpler algorithm. If the target point falls into the scan range, Postgres Pro Shardman performs bidirectional scan starting from this point, advancing at each step in the direction that includes the nearest point. Otherwise, a simple unidirectional scan in the right direction is used.

k-NN search performance has logarithmic dependence on the data size and linear dependence on k. As k tends to the number of rows in the table, the k-NN search time gradually tends to a naive algorithm that consecutively reads the whole table, sorts the data, and returns the first k entries. Since k-NN approach reduces the number of rows that need to be sorted, you can get significant performance benefits if the total number of rows returned by a query is big, or the rows are long. In the worst-case scenario, when all points are equally far from the target point, you have to traverse the whole index, but performance gain can still be achieved because only k entries need to be read from the table, at random.

Another benefit of using k-NN algorithm is simpler logic: you do not need to know the search “radius” in advance, and can resume the search if required.

11.13.1. Examples

Suppose you have the *spots* table containing coordinates of multiple points. To find ten closest points to the given target point, you can run the following query:

```
SELECT coordinates, (coordinates <-> '5.0,5.0'::point) AS dist
FROM spots ORDER BY dist ASC LIMIT 10;
```

coordinates	dist
(3.57192993164062, 6.51727240153665)	2.08362656457647
(3.49502563476562, 6.49134782128243)	2.11874164636854
(3.4393, 6.4473)	2.12848814420001
(3.31787109375, 6.50089913799597)	2.25438592075067
(2.6323, 6.4779)	2.79109148900569
(2.66349792480469, 6.53159856871478)	2.79374947392946
(1.84102535247803, 6.27874198581057)	3.4079762161672
(1.2255, 6.1228)	3.93796014327215
(1.22772216796875, 6.15693947094637)	3.94570513108469


```
(9.6977,4.044503) | 4.79388775494473
(10 rows)
```

If you need to find ten closest points that contain the word "mars" in their names, you can use a multi-column index:

```
CREATE INDEX pt_fts_idx ON spots USING gist(point, to_tsvector('english',asciiname));

SELECT asciiname,point, (point <->'5.0,5.0'::point) AS dist
FROM spots WHERE to_tsvector('english', asciiname)
@@ to_tsquery('english','mars') ORDER BY dist ASC LIMIT 10;
```

Suppose you have the `events` table that describes arbitrary historical events. To select five events that happened close to the launch of the first satellite Sputnik-1 on October 4, 1957, run:

```
SELECT date, event, ('1957-10-04'::date <-> date) AS dist
FROM events ORDER BY dist ASC LIMIT 5;
```

date	event	dist
1957-10-04	Gregory T Linteris, astronaut, is born in New Jersey	0
1957-10-04	USSR launches the first artificial Earth satellite Sputnik-1	0
1957-10-04	"Leave It to Beaver" sitcom debuts on CBS	0
1957-10-03	Lorinc Szabo, Hungarian poet, dies aged 57	1
1957-10-05	All-Stars beat Montreal in 11th NHL All-Star Game	1

(5 rows)

Consider an example of fuzzy search on trigrams. Let's create a trigram index on the `events` table, as follows:

```
CREATE INDEX events_trgm ON events USING gist(event gist_trgm_ops);
```

Now let's find three events that best match the following search query:

```
SELECT date, event, ('george ewashington' <-> event ) AS dist
FROM events ORDER BY dist ASC LIMIT 3;
```

date	event	dist
1732-02-11	George Washington	0.458333
1792-12-05	George Washington is re-elected U.S. President	0.674419
1945-05-12	Jayotis Washington is born	0.714286

(3 rows)

11.13.2. See Also

[Index Types](#)
[Ordering Operators](#)

Chapter 12. Full Text Search

12.1. Introduction

Full Text Searching (or just *text search*) provides the capability to identify natural-language *documents* that satisfy a *query*, and optionally to sort them by relevance to the query. The most common type of search is to find all documents containing given *query terms* and return them in order of their *similarity* to the query. Notions of *query* and *similarity* are very flexible and depend on the specific application. The simplest search considers *query* as a set of words and *similarity* as the frequency of query words in the document.

Textual search operators have existed in databases for years. Postgres Pro has `~`, `~*`, `LIKE`, and `ILIKE` operators for textual data types, but they lack many essential properties required by modern information systems:

- There is no linguistic support, even for English. Regular expressions are not sufficient because they cannot easily handle derived words, e.g., *satisfies* and *satisfy*. You might miss documents that contain *satisfies*, although you probably would like to find them when searching for *satisfy*. It is possible to use `OR` to search for multiple derived forms, but this is tedious and error-prone (some words can have several thousand derivatives).
- They provide no ordering (ranking) of search results, which makes them ineffective when thousands of matching documents are found.
- They tend to be slow because there is no index support, so they must process all documents for every search.

Full text indexing allows documents to be *preprocessed* and an index saved for later rapid searching. Preprocessing includes:

Parsing documents into tokens. It is useful to identify various classes of tokens, e.g., numbers, words, complex words, email addresses, so that they can be processed differently. In principle token classes depend on the specific application, but for most purposes it is adequate to use a predefined set of classes. Postgres Pro uses a *parser* to perform this step. A standard parser is provided, and custom parsers can be created for specific needs.

Converting tokens into lexemes. A lexeme is a string, just like a token, but it has been *normalized* so that different forms of the same word are made alike. For example, normalization almost always includes folding upper-case letters to lower-case, and often involves removal of suffixes (such as *s* or *es* in English). This allows searches to find variant forms of the same word, without tediously entering all the possible variants. Also, this step typically eliminates *stop words*, which are words that are so common that they are useless for searching. (In short, then, tokens are raw fragments of the document text, while lexemes are words that are believed useful for indexing and searching.) Postgres Pro uses *dictionaries* to perform this step. Various standard dictionaries are provided, and custom ones can be created for specific needs.

Storing preprocessed documents optimized for searching. For example, each document can be represented as a sorted array of normalized lexemes. Along with the lexemes it is often desirable to store positional information to use for *proximity ranking*, so that a document that contains a more “dense” region of query words is assigned a higher rank than one with scattered query words.

Dictionaries allow fine-grained control over how tokens are normalized. With appropriate dictionaries, you can:

- Define stop words that should not be indexed.
- Map synonyms to a single word using Ispell.
- Map phrases to a single word using a thesaurus.
- Map different variations of a word to a canonical form using an Ispell dictionary.
- Map different variations of a word to a canonical form using Snowball stemmer rules.

A data type `tsvector` is provided for storing preprocessed documents, along with a type `tsquery` for representing processed queries ([Section 8.11](#)). There are many functions and operators available for

these data types ([Section 9.13](#)), the most important of which is the match operator @@, which we introduce in [Section 12.1.2](#). Full text searches can be accelerated using indexes ([Section 12.9](#)).

12.1.1. What Is a Document?

A *document* is the unit of searching in a full text search system; for example, a magazine article or email message. The text search engine must be able to parse documents and store associations of lexemes (key words) with their parent document. Later, these associations are used to search for documents that contain query words.

For searches within Postgres Pro, a document is normally a textual field within a row of a database table, or possibly a combination (concatenation) of such fields, perhaps stored in several tables or obtained dynamically. In other words, a document can be constructed from different parts for indexing and it might not be stored anywhere as a whole. For example:

```
SELECT title || ' ' || author || ' ' || abstract || ' ' || body AS document
FROM messages
WHERE mid = 12;
```

```
SELECT m.title || ' ' || m.author || ' ' || m.abstract || ' ' || d.body AS document
FROM messages m, docs d
WHERE m.mid = d.did AND m.mid = 12;
```

Note

Actually, in these example queries, `coalesce` should be used to prevent a single `NULL` attribute from causing a `NULL` result for the whole document.

Another possibility is to store the documents as simple text files in the file system. In this case, the database can be used to store the full text index and to execute searches, and some unique identifier can be used to retrieve the document from the file system. However, retrieving files from outside the database requires superuser permissions or special function support, so this is usually less convenient than keeping all the data inside Postgres Pro. Also, keeping everything inside the database allows easy access to document metadata to assist in indexing and display.

For text search purposes, each document must be reduced to the preprocessed `tsvector` format. Searching and ranking are performed entirely on the `tsvector` representation of a document — the original text need only be retrieved when the document has been selected for display to a user. We therefore often speak of the `tsvector` as being the document, but of course it is only a compact representation of the full document.

12.1.2. Basic Text Matching

Full text searching in Postgres Pro is based on the match operator @@, which returns `true` if a `tsvector` (document) matches a `tsquery` (query). It doesn't matter which data type is written first:

```
SELECT 'a fat cat sat on a mat and ate a fat rat'::tsvector @@ 'cat & rat'::tsquery;
?column?
```

```
t
```

```
SELECT 'fat & cow'::tsquery @@ 'a fat cat sat on a mat and ate a fat rat'::tsvector;
?column?
```

```
f
```

As the above example suggests, a `tsquery` is not just raw text, any more than a `tsvector` is. A `tsquery` contains search terms, which must be already-normalized lexemes, and may combine multiple terms using AND, OR, NOT, and FOLLOWED BY operators. (For syntax details see [Section 8.11.2](#).) There are functions `to_tsquery`, `plainto_tsquery`, and `phraseto_tsquery` that are helpful in converting user-

written text into a proper `tsquery`, primarily by normalizing words appearing in the text. Similarly, `to_tsvector` is used to parse and normalize a document string. So in practice a text search match would look more like this:

```
SELECT to_tsvector('fat cats ate fat rats') @@ to_tsquery('fat & rat');
?column?
-----
t
```

Observe that this match would not succeed if written as

```
SELECT 'fat cats ate fat rats'::tsvector @@ to_tsquery('fat & rat');
?column?
-----
f
```

since here no normalization of the word `rats` will occur. The elements of a `tsvector` are lexemes, which are assumed already normalized, so `rats` does not match `rat`.

The `@@` operator also supports `text` input, allowing explicit conversion of a text string to `tsvector` or `tsquery` to be skipped in simple cases. The variants available are:

```
tsvector @@ tsquery
tsquery  @@ tsvector
text     @@ tsquery
text     @@ text
```

The first two of these we saw already. The form `text @@ tsquery` is equivalent to `to_tsvector(x) @@ y`. The form `text @@ text` is equivalent to `to_tsvector(x) @@ plainto_tsquery(y)`.

Within a `tsquery`, the `&` (AND) operator specifies that both its arguments must appear in the document to have a match. Similarly, the `|` (OR) operator specifies that at least one of its arguments must appear, while the `!` (NOT) operator specifies that its argument must *not* appear in order to have a match. For example, the query `fat & ! rat` matches documents that contain `fat` but not `rat`.

Searching for phrases is possible with the help of the `<->` (FOLLOWED BY) `tsquery` operator, which matches only if its arguments have matches that are adjacent and in the given order. For example:

```
SELECT to_tsvector('fatal error') @@ to_tsquery('fatal <-> error');
?column?
-----
t

SELECT to_tsvector('error is not fatal') @@ to_tsquery('fatal <-> error');
?column?
-----
f
```

There is a more general version of the FOLLOWED BY operator having the form `<N>`, where `N` is an integer standing for the difference between the positions of the matching lexemes. `<1>` is the same as `<->`, while `<2>` allows exactly one other lexeme to appear between the matches, and so on. The `phraseto_tsquery` function makes use of this operator to construct a `tsquery` that can match a multi-word phrase when some of the words are stop words. For example:

```
SELECT phraseto_tsquery('cats ate rats');
phraseto_tsquery
-----
'cat' <-> 'ate' <-> 'rat'

SELECT phraseto_tsquery('the cats ate the rats');
phraseto_tsquery
-----
'cat' <-> 'ate' <2> 'rat'
```

A special case that's sometimes useful is that `<0>` can be used to require that two patterns match the same word.

Parentheses can be used to control nesting of the `tsquery` operators. Without parentheses, `|` binds least tightly, then `&`, then `<->`, and `!` most tightly.

It's worth noticing that the AND/OR/NOT operators mean something subtly different when they are within the arguments of a FOLLOWED BY operator than when they are not, because within FOLLOWED BY the exact position of the match is significant. For example, normally `!x` matches only documents that do not contain `x` anywhere. But `!x <-> y` matches `y` if it is not immediately after an `x`; an occurrence of `x` elsewhere in the document does not prevent a match. Another example is that `x & y` normally only requires that `x` and `y` both appear somewhere in the document, but `(x & y) <-> z` requires `x` and `y` to match at the same place, immediately before a `z`. Thus this query behaves differently from `x <-> z & y <-> z`, which will match a document containing two separate sequences `x z` and `y z`. (This specific query is useless as written, since `x` and `y` could not match at the same place; but with more complex situations such as prefix-match patterns, a query of this form could be useful.)

12.1.3. Configurations

The above are all simple text search examples. As mentioned before, full text search functionality includes the ability to do many more things: skip indexing certain words (stop words), process synonyms, and use sophisticated parsing, e.g., parse based on more than just white space. This functionality is controlled by *text search configurations*. Postgres Pro comes with predefined configurations for many languages, and you can easily create your own configurations. (psql's `\df` command shows all available configurations.)

During installation an appropriate configuration is selected and [default_text_search_config](#) is set accordingly in `postgresql.conf`. If you are using the same text search configuration for the entire cluster you can use the value in `postgresql.conf`. To use different configurations throughout the cluster but the same configuration within any one database, use `ALTER DATABASE ... SET`. Otherwise, you can set `default_text_search_config` in each session.

Each text search function that depends on a configuration has an optional `regconfig` argument, so that the configuration to use can be specified explicitly. `default_text_search_config` is used only when this argument is omitted.

To make it easier to build custom text search configurations, a configuration is built up from simpler database objects. Postgres Pro's text search facility provides four types of configuration-related database objects:

- *Text search parsers* break documents into tokens and classify each token (for example, as words or numbers).
- *Text search dictionaries* convert tokens to normalized form and reject stop words.
- *Text search templates* provide the functions underlying dictionaries. (A dictionary simply specifies a template and a set of parameters for the template.)
- *Text search configurations* select a parser and a set of dictionaries to use to normalize the tokens produced by the parser.

Text search parsers and templates are built from low-level C functions; therefore it requires C programming ability to develop new ones, and superuser privileges to install one into a database. (There are examples of add-on parsers and templates in the `contrib/` area of the Postgres Pro distribution.) Since dictionaries and configurations just parameterize and connect together some underlying parsers and templates, no special privilege is needed to create a new dictionary or configuration. Examples of creating custom dictionaries and configurations appear later in this chapter.

12.2. Tables and Indexes

The examples in the previous section illustrated full text matching using simple constant strings. This section shows how to search table data, optionally using indexes.

12.2.1. Searching a Table

It is possible to do a full text search without an index. A simple query to print the `title` of each row that contains the word `friend` in its `body` field is:

```
SELECT title
FROM pgweb
WHERE to_tsvector('english', body) @@ to_tsquery('english', 'friend');
```

This will also find related words such as `friends` and `friendly`, since all these are reduced to the same normalized lexeme.

The query above specifies that the `english` configuration is to be used to parse and normalize the strings. Alternatively we could omit the configuration parameters:

```
SELECT title
FROM pgweb
WHERE to_tsvector(body) @@ to_tsquery('friend');
```

This query will use the configuration set by [default_text_search_config](#).

A more complex example is to select the ten most recent documents that contain `create` and `table` in the `title` or `body`:

```
SELECT title
FROM pgweb
WHERE to_tsvector(title || ' ' || body) @@ to_tsquery('create & table')
ORDER BY last_mod_date DESC
LIMIT 10;
```

For clarity we omitted the `coalesce` function calls which would be needed to find rows that contain `NULL` in one of the two fields.

Although these queries will work without an index, most applications will find this approach too slow, except perhaps for occasional ad-hoc searches. Practical use of text searching usually requires creating an index.

12.2.2. Creating Indexes

We can create a GIN index ([Section 12.9](#)) to speed up text searches:

```
CREATE INDEX pgweb_idx ON pgweb USING GIN (to_tsvector('english', body));
```

Notice that the 2-argument version of `to_tsvector` is used. Only text search functions that specify a configuration name can be used in expression indexes ([Section 11.7](#)). This is because the index contents must be unaffected by [default_text_search_config](#). If they were affected, the index contents might be inconsistent because different entries could contain `tsvectors` that were created with different text search configurations, and there would be no way to guess which was which. It would be impossible to dump and restore such an index correctly.

Because the two-argument version of `to_tsvector` was used in the index above, only a query reference that uses the 2-argument version of `to_tsvector` with the same configuration name will use that index. That is, `WHERE to_tsvector('english', body) @@ 'a & b'` can use the index, but `WHERE to_tsvector(body) @@ 'a & b'` cannot. This ensures that an index will be used only with the same configuration used to create the index entries.

It is possible to set up more complex expression indexes wherein the configuration name is specified by another column, e.g.:

```
CREATE INDEX pgweb_idx ON pgweb USING GIN (to_tsvector(config_name, body));
```

where `config_name` is a column in the `pgweb` table. This allows mixed configurations in the same index while recording which configuration was used for each index entry. This would be useful, for example, if the document collection contained documents in different languages. Again, queries that are meant to use the index must be phrased to match, e.g., `WHERE to_tsvector(config_name, body) @@ 'a & b'`.

Indexes can even concatenate columns:

```
CREATE INDEX pgweb_idx ON pgweb USING GIN (to_tsvector('english', title || ' ' ||
body));
```

Another approach is to create a separate `tsvector` column to hold the output of `to_tsvector`. To keep this column automatically up to date with its source data, use a stored generated column. This example is a concatenation of `title` and `body`, using `coalesce` to ensure that one field will still be indexed when the other is `NULL`:

```
ALTER TABLE pgweb
  ADD COLUMN textsearchable_index_col tsvector
  GENERATED ALWAYS AS (to_tsvector('english', coalesce(title, '') || ' '
|| coalesce(body, ''))) STORED;
```

Then we create a GIN index to speed up the search:

```
CREATE INDEX textsearch_idx ON pgweb USING GIN (textsearchable_index_col);
```

Now we are ready to perform a fast full text search:

```
SELECT title
FROM pgweb
WHERE textsearchable_index_col @@ to_tsquery('create & table')
ORDER BY last_mod_date DESC
LIMIT 10;
```

One advantage of the separate-column approach over an expression index is that it is not necessary to explicitly specify the text search configuration in queries in order to make use of the index. As shown in the example above, the query can depend on `default_text_search_config`. Another advantage is that searches will be faster, since it will not be necessary to redo the `to_tsvector` calls to verify index matches. (This is more important when using a GiST index than a GIN index; see [Section 12.9](#).) The expression-index approach is simpler to set up, however, and it requires less disk space since the `tsvector` representation is not stored explicitly.

12.3. Controlling Text Search

To implement full text searching there must be a function to create a `tsvector` from a document and a `tsquery` from a user query. Also, we need to return results in a useful order, so we need a function that compares documents with respect to their relevance to the query. It's also important to be able to display the results nicely. Postgres Pro provides support for all of these functions.

12.3.1. Parsing Documents

Postgres Pro provides the function `to_tsvector` for converting a document to the `tsvector` data type.

`to_tsvector([config regconfig,] document text)` returns `tsvector`

`to_tsvector` parses a textual document into tokens, reduces the tokens to lexemes, and returns a `tsvector` which lists the lexemes together with their positions in the document. The document is processed according to the specified or default text search configuration. Here is a simple example:

```
SELECT to_tsvector('english', 'a fat  cat sat on a mat - it ate a fat rats');
       to_tsvector
-----
'ate':9 'cat':3 'fat':2,11 'mat':7 'rat':12 'sat':4
```

In the example above we see that the resulting `tsvector` does not contain the words `a`, `on`, or `it`, the word `rats` became `rat`, and the punctuation sign `-` was ignored.

The `to_tsvector` function internally calls a parser which breaks the document text into tokens and assigns a type to each token. For each token, a list of dictionaries ([Section 12.6](#)) is consulted, where

the list can vary depending on the token type. The first dictionary that *recognizes* the token emits one or more normalized *lexemes* to represent the token. For example, `rats` became `rat` because one of the dictionaries recognized that the word `rats` is a plural form of `rat`. Some words are recognized as *stop words* ([Section 12.6.1](#)), which causes them to be ignored since they occur too frequently to be useful in searching. In our example these are `a`, `on`, and `it`. If no dictionary in the list recognizes the token then it is also ignored. In this example that happened to the punctuation sign – because there are in fact no dictionaries assigned for its token type (`Space symbols`), meaning space tokens will never be indexed. The choices of parser, dictionaries and which types of tokens to index are determined by the selected text search configuration ([Section 12.7](#)). It is possible to have many different configurations in the same database, and predefined configurations are available for various languages. In our example we used the default configuration `english` for the English language.

The function `setweight` can be used to label the entries of a `tsvector` with a given *weight*, where a weight is one of the letters A, B, C, or D. This is typically used to mark entries coming from different parts of a document, such as title versus body. Later, this information can be used for ranking of search results.

Because `to_tsvector(NULL)` will return `NULL`, it is recommended to use `coalesce` whenever a field might be null. Here is the recommended method for creating a `tsvector` from a structured document:

```
UPDATE tt SET ti =
    setweight(to_tsvector(coalesce(title,'')), 'A')      ||
    setweight(to_tsvector(coalesce(keyword,'')), 'B')   ||
    setweight(to_tsvector(coalesce(abstract,'')), 'C')  ||
    setweight(to_tsvector(coalesce(body,'')), 'D');
```

Here we have used `setweight` to label the source of each lexeme in the finished `tsvector`, and then merged the labeled `tsvector` values using the `tsvector` concatenation operator `||`. ([Section 12.4.1](#) gives details about these operations.)

12.3.2. Parsing Queries

Postgres Pro provides the functions `to_tsquery`, `plainto_tsquery`, `phraseto_tsquery` and `web-search_to_tsquery` for converting a query to the `tsquery` data type. `to_tsquery` offers access to more features than either `plainto_tsquery` or `phraseto_tsquery`, but it is less forgiving about its input. `web-search_to_tsquery` is a simplified version of `to_tsquery` with an alternative syntax, similar to the one used by web search engines.

`to_tsquery([config regconfig,] querytext text)` returns `tsquery`

`to_tsquery` creates a `tsquery` value from `querytext`, which must consist of single tokens separated by the `tsquery` operators `&` (AND), `|` (OR), `!` (NOT), and `<->` (FOLLOWED BY), possibly grouped using parentheses. In other words, the input to `to_tsquery` must already follow the general rules for `tsquery` input, as described in [Section 8.11.2](#). The difference is that while basic `tsquery` input takes the tokens at face value, `to_tsquery` normalizes each token into a lexeme using the specified or default configuration, and discards any tokens that are stop words according to the configuration. For example:

```
SELECT to_tsquery('english', 'The & Fat & Rats');
 to_tsquery
-----
'fat' & 'rat'
```

As in basic `tsquery` input, `weight(s)` can be attached to each lexeme to restrict it to match only `tsvector` lexemes of those `weight(s)`. For example:

```
SELECT to_tsquery('english', 'Fat | Rats:AB');
 to_tsquery
-----
'fat' | 'rat':AB
```

Also, `*` can be attached to a lexeme to specify prefix matching:

```
SELECT to_tsquery('supern:*A & star:A*B');
```



```
to_tsquery
```

```
-----  
'supern':*A & 'star':*AB
```

Such a lexeme will match any word in a `tsvector` that begins with the given string.

`to_tsquery` can also accept single-quoted phrases. This is primarily useful when the configuration includes a thesaurus dictionary that may trigger on such phrases. In the example below, a thesaurus contains the rule `supernovae stars : sn`:

```
SELECT to_tsquery('supernovae stars' & !crab');  
to_tsquery  
-----  
'sn' & !'crab'
```

Without quotes, `to_tsquery` will generate a syntax error for tokens that are not separated by an AND, OR, or FOLLOWED BY operator.

`plainto_tsquery([config regconfig,] querytext text)` returns `tsquery`

`plainto_tsquery` transforms the unformatted text `querytext` to a `tsquery` value. The text is parsed and normalized much as for `to_tsvector`, then the `&` (AND) `tsquery` operator is inserted between surviving words.

Example:

```
SELECT plainto_tsquery('english', 'The Fat Rats');  
plainto_tsquery  
-----  
'fat' & 'rat'
```

Note that `plainto_tsquery` will not recognize `tsquery` operators, weight labels, or prefix-match labels in its input:

```
SELECT plainto_tsquery('english', 'The Fat & Rats:C');  
plainto_tsquery  
-----  
'fat' & 'rat' & 'c'
```

Here, all the input punctuation was discarded.

`phraseto_tsquery([config regconfig,] querytext text)` returns `tsquery`

`phraseto_tsquery` behaves much like `plainto_tsquery`, except that it inserts the `<->` (FOLLOWED BY) operator between surviving words instead of the `&` (AND) operator. Also, stop words are not simply discarded, but are accounted for by inserting `<N>` operators rather than `<->` operators. This function is useful when searching for exact lexeme sequences, since the FOLLOWED BY operators check lexeme order not just the presence of all the lexemes.

Example:

```
SELECT phraseto_tsquery('english', 'The Fat Rats');  
phraseto_tsquery  
-----  
'fat' <-> 'rat'
```

Like `plainto_tsquery`, the `phraseto_tsquery` function will not recognize `tsquery` operators, weight labels, or prefix-match labels in its input:

```
SELECT phraseto_tsquery('english', 'The Fat & Rats:C');  
phraseto_tsquery  
-----  
'fat' <-> 'rat' <-> 'c'
```

`websearch_to_tsquery([config regconfig,] querytext text)` returns `tsquery`

`websearch_to_tsquery` creates a `tsquery` value from `querytext` using an alternative syntax in which simple unformatted text is a valid query. Unlike `plainto_tsquery` and `phraseto_tsquery`, it also recognizes certain operators. Moreover, this function will never raise syntax errors, which makes it possible to use raw user-supplied input for search. The following syntax is supported:

- `unquoted text`: text not inside quote marks will be converted to terms separated by `&` operators, as if processed by `plainto_tsquery`.
- `"quoted text"`: text inside quote marks will be converted to terms separated by `<->` operators, as if processed by `phraseto_tsquery`.
- `OR`: the word “or” will be converted to the `|` operator.
- `-`: a dash will be converted to the `!` operator.

Other punctuation is ignored. So like `plainto_tsquery` and `phraseto_tsquery`, the `websearch_to_tsquery` function will not recognize `tsquery` operators, weight labels, or prefix-match labels in its input.

Examples:

```
SELECT websearch_to_tsquery('english', 'The fat rats');
       websearch_to_tsquery
```

```
-----
'fat' & 'rat'
(1 row)
```

```
SELECT websearch_to_tsquery('english', '"supernovae stars" -crab');
       websearch_to_tsquery
```

```
-----
'supernova' <-> 'star' & !'crab'
(1 row)
```

```
SELECT websearch_to_tsquery('english', '"sad cat" or "fat rat"');
       websearch_to_tsquery
```

```
-----
'sad' <-> 'cat' | 'fat' <-> 'rat'
(1 row)
```

```
SELECT websearch_to_tsquery('english', 'signal -"segmentation fault"');
       websearch_to_tsquery
```

```
-----
'signal' & !( 'segment' <-> 'fault' )
(1 row)
```

```
SELECT websearch_to_tsquery('english', '""')(dummy \ query <->');
       websearch_to_tsquery
```

```
-----
'dummi' & 'queri'
(1 row)
```

12.3.3. Ranking Search Results

Ranking attempts to measure how relevant documents are to a particular query, so that when there are many matches the most relevant ones can be shown first. Postgres Pro provides two predefined ranking functions, which take into account lexical, proximity, and structural information; that is, they consider how often the query terms appear in the document, how close together the terms are in the document, and how important is the part of the document where they occur. However, the concept of relevancy is vague and very application-specific. Different applications might require additional information for ranking, e.g., document modification time. The built-in ranking functions are only examples. You can write your own ranking functions and/or combine their results with additional factors to fit your specific needs.

The two ranking functions currently available are:

```
ts_rank([ weights float4[], ] vector tsvector, query tsquery [, normalization integer
]) returns float4
```

Ranks vectors based on the frequency of their matching lexemes.

```
ts_rank_cd([ weights float4[], ] vector tsvector, query tsquery [, normalization integer
]) returns float4
```

This function computes the *cover density* ranking for the given document vector and query, as described in Clarke, Cormack, and Tudhope's "Relevance Ranking for One to Three Term Queries" in the journal "Information Processing and Management", 1999. Cover density is similar to `ts_rank` ranking except that the proximity of matching lexemes to each other is taken into consideration.

This function requires lexeme positional information to perform its calculation. Therefore, it ignores any “stripped” lexemes in the `tsvector`. If there are no unstripped lexemes in the input, the result will be zero. (See [Section 12.4.1](#) for more information about the `strip` function and positional information in `tsvectors`.)

For both these functions, the optional `weights` argument offers the ability to weigh word instances more or less heavily depending on how they are labeled. The weight arrays specify how heavily to weigh each category of word, in the order:

```
{D-weight, C-weight, B-weight, A-weight}
```

If no `weights` are provided, then these defaults are used:

```
{0.1, 0.2, 0.4, 1.0}
```

Typically weights are used to mark words from special areas of the document, like the title or an initial abstract, so they can be treated with more or less importance than words in the document body.

Since a longer document has a greater chance of containing a query term it is reasonable to take into account document size, e.g., a hundred-word document with five instances of a search word is probably more relevant than a thousand-word document with five instances. Both ranking functions take an integer `normalization` option that specifies whether and how a document's length should impact its rank. The integer option controls several behaviors, so it is a bit mask: you can specify one or more behaviors using `|` (for example, `2|4`).

- 0 (the default) ignores the document length
- 1 divides the rank by 1 + the logarithm of the document length
- 2 divides the rank by the document length
- 4 divides the rank by the mean harmonic distance between extents (this is implemented only by `ts_rank_cd`)
- 8 divides the rank by the number of unique words in document
- 16 divides the rank by 1 + the logarithm of the number of unique words in document
- 32 divides the rank by itself + 1

If more than one flag bit is specified, the transformations are applied in the order listed.

It is important to note that the ranking functions do not use any global information, so it is impossible to produce a fair normalization to 1% or 100% as sometimes desired. Normalization option 32 (`rank/(rank+1)`) can be applied to scale all ranks into the range zero to one, but of course this is just a cosmetic change; it will not affect the ordering of the search results.

Here is an example that selects only the ten highest-ranked matches:

```
SELECT title, ts_rank_cd(textsearch, query) AS rank
FROM apod, to_tsquery('neutrino|(dark & matter)') query
WHERE query @@ textsearch
ORDER BY rank DESC
LIMIT 10;
```

title	rank
-------	------

Neutrinos in the Sun	3.1
The Sudbury Neutrino Detector	2.4
A MACHO View of Galactic Dark Matter	2.01317
Hot Gas and Dark Matter	1.91171
The Virgo Cluster: Hot Plasma and Dark Matter	1.90953
Rafting for Solar Neutrinos	1.9
NGC 4650A: Strange Galaxy and Dark Matter	1.85774
Hot Gas and Dark Matter	1.6123
Ice Fishing for Cosmic Neutrinos	1.6
Weak Lensing Distorts the Universe	0.818218

This is the same example using normalized ranking:

```
SELECT title, ts_rank_cd(textsearch, query, 32 /* rank/(rank+1) */ ) AS rank
FROM apod, to_tsquery('neutrino|(dark & matter)') query
WHERE query @@ textsearch
ORDER BY rank DESC
LIMIT 10;
```

Neutrinos in the Sun	0.756097569485493
The Sudbury Neutrino Detector	0.705882361190954
A MACHO View of Galactic Dark Matter	0.668123210574724
Hot Gas and Dark Matter	0.65655958650282
The Virgo Cluster: Hot Plasma and Dark Matter	0.656301290640973
Rafting for Solar Neutrinos	0.655172410958162
NGC 4650A: Strange Galaxy and Dark Matter	0.650072921219637
Hot Gas and Dark Matter	0.617195790024749
Ice Fishing for Cosmic Neutrinos	0.615384618911517
Weak Lensing Distorts the Universe	0.450010798361481

Ranking can be expensive since it requires consulting the `tsvector` of each matching document, which can be I/O bound and therefore slow. Unfortunately, it is almost impossible to avoid since practical queries often result in large numbers of matches.

12.3.4. Highlighting Results

To present search results it is ideal to show a part of each document and how it is related to the query. Usually, search engines show fragments of the document with marked search terms. Postgres Pro provides a function `ts_headline` that implements this functionality.

```
ts_headline([ config regconfig, ] document text, query tsquery [, options text ])
returns text
```

`ts_headline` accepts a document along with a query, and returns an excerpt from the document in which terms from the query are highlighted. Specifically, the function will use the query to select relevant text fragments, and then highlight all words that appear in the query, even if those word positions do not match the query's restrictions. The configuration to be used to parse the document can be specified by `config`; if `config` is omitted, the `default_text_search_config` configuration is used.

If an `options` string is specified it must consist of a comma-separated list of one or more `option=value` pairs. The available options are:

- `MaxWords`, `MinWords` (integers): these numbers determine the longest and shortest headlines to output. The default values are 35 and 15.
- `ShortWord` (integer): words of this length or less will be dropped at the start and end of a headline, unless they are query terms. The default value of three eliminates common English articles.
- `HighlightAll` (boolean): if `true` the whole document will be used as the headline, ignoring the preceding three parameters. The default is `false`.

- `MaxFragments` (integer): maximum number of text fragments to display. The default value of zero selects a non-fragment-based headline generation method. A value greater than zero selects fragment-based headline generation (see below).
- `StartSel`, `StopSel` (strings): the strings with which to delimit query words appearing in the document, to distinguish them from other excerpted words. The default values are “” and “”, which can be suitable for HTML output (but see the warning below).
- `FragmentDelimiter` (string): When more than one fragment is displayed, the fragments will be separated by this string. The default is “ ... ”.

Warning: Cross-site scripting (XSS) safety

The output from `ts_headline` is not guaranteed to be safe for direct inclusion in web pages. When `HighlightAll` is `false` (the default), some simple XML tags are removed from the document, but this is not guaranteed to remove all HTML markup. Therefore, this does not provide an effective defense against attacks such as cross-site scripting (XSS) attacks, when working with untrusted input. To guard against such attacks, all HTML markup should be removed from the input document, or an HTML sanitizer should be used on the output.

These option names are recognized case-insensitively. You must double-quote string values if they contain spaces or commas.

In non-fragment-based headline generation, `ts_headline` locates matches for the given *query* and chooses a single one to display, preferring matches that have more query words within the allowed headline length. In fragment-based headline generation, `ts_headline` locates the query matches and splits each match into “fragments” of no more than `MaxWords` words each, preferring fragments with more query words, and when possible “stretching” fragments to include surrounding words. The fragment-based mode is thus more useful when the query matches span large sections of the document, or when it's desirable to display multiple matches. In either mode, if no query matches can be identified, then a single fragment of the first `MinWords` words in the document will be displayed.

For example:

```
SELECT ts_headline('english',
  'The most common type of search
is to find all documents containing given query terms
and return them in order of their similarity to the
query.',
  to_tsquery('english', 'query & similarity'));
      ts_headline
-----
containing given <b>query</b> terms                +
and return them in order of their <b>similarity</b> to the+
<b>query</b>.
```

```
SELECT ts_headline('english',
  'Search terms may occur
many times in a document,
requiring ranking of the search matches to decide which
occurrences to display in the result.',
  to_tsquery('english', 'search & term'),
  'MaxFragments=10, MaxWords=7, MinWords=3, StartSel=<<, StopSel=>>');
      ts_headline
-----
<<Search>> <<terms>> may occur                    +
many times ... ranking of the <<search>> matches to decide
```

`ts_headline` uses the original document, not a `tsvector` summary, so it can be slow and should be used with care.

12.4. Additional Features

This section describes additional functions and operators that are useful in connection with text search.

12.4.1. Manipulating Documents

[Section 12.3.1](#) showed how raw textual documents can be converted into `tsvector` values. Postgres Pro also provides functions and operators that can be used to manipulate documents that are already in `tsvector` form.

```
tsvector || tsvector
```

The `tsvector` concatenation operator returns a vector which combines the lexemes and positional information of the two vectors given as arguments. Positions and weight labels are retained during the concatenation. Positions appearing in the right-hand vector are offset by the largest position mentioned in the left-hand vector, so that the result is nearly equivalent to the result of performing `to_tsvector` on the concatenation of the two original document strings. (The equivalence is not exact, because any stop-words removed from the end of the left-hand argument will not affect the result, whereas they would have affected the positions of the lexemes in the right-hand argument if textual concatenation were used.)

One advantage of using concatenation in the vector form, rather than concatenating text before applying `to_tsvector`, is that you can use different configurations to parse different sections of the document. Also, because the `setweight` function marks all lexemes of the given vector the same way, it is necessary to parse the text and do `setweight` before concatenating if you want to label different parts of the document with different weights.

```
setweight(vector tsvector, weight "char") returns tsvector
```

`setweight` returns a copy of the input vector in which every position has been labeled with the given *weight*, either A, B, C, or D. (D is the default for new vectors and as such is not displayed on output.) These labels are retained when vectors are concatenated, allowing words from different parts of a document to be weighted differently by ranking functions.

Note that weight labels apply to *positions*, not *lexemes*. If the input vector has been stripped of positions then `setweight` does nothing.

```
length(vector tsvector) returns integer
```

Returns the number of lexemes stored in the vector.

```
strip(vector tsvector) returns tsvector
```

Returns a vector that lists the same lexemes as the given vector, but lacks any position or weight information. The result is usually much smaller than an unstripped vector, but it is also less useful. Relevance ranking does not work as well on stripped vectors as unstripped ones. Also, the `<->` (FOLLOWED BY) `tsquery` operator will never match stripped input, since it cannot determine the distance between lexeme occurrences.

A full list of `tsvector`-related functions is available in [Table 9.43](#).

12.4.2. Manipulating Queries

[Section 12.3.2](#) showed how raw textual queries can be converted into `tsquery` values. Postgres Pro also provides functions and operators that can be used to manipulate queries that are already in `tsquery` form.

```
tsquery && tsquery
```

Returns the AND-combination of the two given queries.

```
tsquery || tsquery
```

Returns the OR-combination of the two given queries.

!! tsquery

Returns the negation (NOT) of the given query.

tsquery <-> tsquery

Returns a query that searches for a match to the first given query immediately followed by a match to the second given query, using the <-> (FOLLOWED BY) *tsquery* operator. For example:

```
SELECT to_tsquery('fat') <-> to_tsquery('cat | rat');
       ?column?
-----
'fat' <-> ( 'cat' | 'rat' )
```

tsquery_phrase(query1 tsquery, query2 tsquery [, distance integer]) returns *tsquery*

Returns a query that searches for a match to the first given query followed by a match to the second given query at a distance of exactly *distance* lexemes, using the <N> *tsquery* operator. For example:

```
SELECT tsquery_phrase(to_tsquery('fat'), to_tsquery('cat'), 10);
       tsquery_phrase
-----
'fat' <10> 'cat'
```

numnode(query tsquery) returns *integer*

Returns the number of nodes (lexemes plus operators) in a *tsquery*. This function is useful to determine if the *query* is meaningful (returns > 0), or contains only stop words (returns 0). Examples:

```
SELECT numnode(plainto_tsquery('the any'));
NOTICE: query contains only stopword(s) or doesn't contain lexeme(s), ignored
       numnode
-----
           0

SELECT numnode('foo & bar'::tsquery);
       numnode
-----
           3
```

querytree(query tsquery) returns *text*

Returns the portion of a *tsquery* that can be used for searching an index. This function is useful for detecting unindexable queries, for example those containing only stop words or only negated terms. For example:

```
SELECT querytree(to_tsquery('defined'));
       querytree
-----
'defin'

SELECT querytree(to_tsquery('!defined'));
       querytree
-----
T
```

12.4.2.1. Query Rewriting

The *ts_rewrite* family of functions search a given *tsquery* for occurrences of a target subquery, and replace each occurrence with a substitute subquery. In essence this operation is a *tsquery*-specific version of substring replacement. A target and substitute combination can be thought of as a *query rewrite rule*. A collection of such rewrite rules can be a powerful search aid. For example, you can expand the search using synonyms (e.g., *new york*, *big apple*, *nyc*, *gotham*) or narrow the search

to direct the user to some hot topic. There is some overlap in functionality between this feature and thesaurus dictionaries ([Section 12.6.4](#)). However, you can modify a set of rewrite rules on-the-fly without reindexing, whereas updating a thesaurus requires reindexing to be effective.

`ts_rewrite (query tsquery, target tsquery, substitute tsquery)` returns `tsquery`

This form of `ts_rewrite` simply applies a single rewrite rule: *target* is replaced by *substitute* wherever it appears in *query*. For example:

```
SELECT ts_rewrite('a & b'::tsquery, 'a'::tsquery, 'c'::tsquery);
      ts_rewrite
-----
      'b' & 'c'
```

`ts_rewrite (query tsquery, select text)` returns `tsquery`

This form of `ts_rewrite` accepts a starting *query* and an SQL *select* command, which is given as a text string. The *select* must yield two columns of `tsquery` type. For each row of the *select* result, occurrences of the first column value (the target) are replaced by the second column value (the substitute) within the current *query* value. For example:

```
CREATE TABLE aliases (t tsquery PRIMARY KEY, s tsquery);
INSERT INTO aliases VALUES('a', 'c');

SELECT ts_rewrite('a & b'::tsquery, 'SELECT t,s FROM aliases');
      ts_rewrite
-----
      'b' & 'c'
```

Note that when multiple rewrite rules are applied in this way, the order of application can be important; so in practice you will want the source query to `ORDER BY` some ordering key.

Let's consider a real-life astronomical example. We'll expand query `supernovae` using table-driven rewriting rules:

```
CREATE TABLE aliases (t tsquery primary key, s tsquery);
INSERT INTO aliases VALUES(to_tsquery('supernovae'), to_tsquery('supernovae|sn'));

SELECT ts_rewrite(to_tsquery('supernovae & crab'), 'SELECT * FROM aliases');
      ts_rewrite
-----
      'crab' & ( 'supernova' | 'sn' )
```

We can change the rewriting rules just by updating the table:

```
UPDATE aliases
SET s = to_tsquery('supernovae|sn & !nebulae')
WHERE t = to_tsquery('supernovae');

SELECT ts_rewrite(to_tsquery('supernovae & crab'), 'SELECT * FROM aliases');
      ts_rewrite
-----
      'crab' & ( 'supernova' | 'sn' & '!nebula' )
```

Rewriting can be slow when there are many rewriting rules, since it checks every rule for a possible match. To filter out obvious non-candidate rules we can use the containment operators for the `tsquery` type. In the example below, we select only those rules which might match the original query:

```
SELECT ts_rewrite('a & b'::tsquery,
                  'SELECT t,s FROM aliases WHERE 'a & b'::tsquery @> t');
      ts_rewrite
-----
      'b' & 'c'
```


12.4.3. Triggers for Automatic Updates

Note

The method described in this section has been obsoleted by the use of stored generated columns, as described in [Section 12.2.2](#).

When using a separate column to store the `tsvector` representation of your documents, it is necessary to create a trigger to update the `tsvector` column when the document content columns change. Two built-in trigger functions are available for this, or you can write your own.

```
tsvector_update_trigger(tsvector_column_name, config_name, text_column_name [, ... ])
tsvector_update_trigger_column(tsvector_column_name,
                               config_column_name, text_column_name [, ... ])
```

These trigger functions automatically compute a `tsvector` column from one or more textual columns, under the control of parameters specified in the `CREATE TRIGGER` command. An example of their use is:

```
CREATE TABLE messages (
    title      text,
    body       text,
    tsv        tsvector
);

CREATE TRIGGER tsvectorupdate BEFORE INSERT OR UPDATE
ON messages FOR EACH ROW EXECUTE FUNCTION
tsvector_update_trigger(tsv, 'pg_catalog.english', title, body);

INSERT INTO messages VALUES('title here', 'the body text is here');

SELECT * FROM messages;
   title      |          body          |          tsv
-----+-----+-----
title here | the body text is here | 'bodi':4 'text':5 'titl':1

SELECT title, body FROM messages WHERE tsv @@ to_tsquery('title & body');
   title      |          body
-----+-----
title here | the body text is here
```

Having created this trigger, any change in `title` or `body` will automatically be reflected into `tsv`, without the application having to worry about it.

The first trigger argument must be the name of the `tsvector` column to be updated. The second argument specifies the text search configuration to be used to perform the conversion. For `tsvector_update_trigger`, the configuration name is simply given as the second trigger argument. It must be schema-qualified as shown above, so that the trigger behavior will not change with changes in `search_path`. For `tsvector_update_trigger_column`, the second trigger argument is the name of another table column, which must be of type `regconfig`. This allows a per-row selection of configuration to be made. The remaining argument(s) are the names of textual columns (of type `text`, `varchar`, or `char`). These will be included in the document in the order given. NULL values will be skipped (but the other columns will still be indexed).

A limitation of these built-in triggers is that they treat all the input columns alike. To process columns differently — for example, to weight `title` differently from `body` — it is necessary to write a custom trigger. Here is an example using PL/pgSQL as the trigger language:

```
CREATE FUNCTION messages_trigger() RETURNS trigger AS $$
begin
```

```

new.tsv :=
    setweight(to_tsvector('pg_catalog.english', coalesce(new.title, '')), 'A') ||
    setweight(to_tsvector('pg_catalog.english', coalesce(new.body, '')), 'D');
return new;
end
$$ LANGUAGE plpgsql;

```

```

CREATE TRIGGER tsvectorupdate BEFORE INSERT OR UPDATE
    ON messages FOR EACH ROW EXECUTE FUNCTION messages_trigger();

```

Keep in mind that it is important to specify the configuration name explicitly when creating `tsvector` values inside triggers, so that the column's contents will not be affected by changes to `default_text_search_config`. Failure to do this is likely to lead to problems such as search results changing after a dump and restore.

12.4.4. Gathering Document Statistics

The function `ts_stat` is useful for checking your configuration and for finding stop-word candidates.

```

ts_stat(sqlquery text, [ weights text, ]
        OUT word text, OUT ndoc integer,
        OUT nentry integer) returns setof record

```

`sqlquery` is a text value containing an SQL query which must return a single `tsvector` column. `ts_stat` executes the query and returns statistics about each distinct lexeme (word) contained in the `tsvector` data. The columns returned are

- `word text` — the value of a lexeme
- `ndoc integer` — number of documents (`tsvectors`) the word occurred in
- `nentry integer` — total number of occurrences of the word

If `weights` is supplied, only occurrences having one of those weights are counted.

For example, to find the ten most frequent words in a document collection:

```

SELECT * FROM ts_stat('SELECT vector FROM apod')
ORDER BY nentry DESC, ndoc DESC, word
LIMIT 10;

```

The same, but counting only word occurrences with weight A or B:

```

SELECT * FROM ts_stat('SELECT vector FROM apod', 'ab')
ORDER BY nentry DESC, ndoc DESC, word
LIMIT 10;

```

12.5. Parsers

Text search parsers are responsible for splitting raw document text into *tokens* and identifying each token's type, where the set of possible types is defined by the parser itself. Note that a parser does not modify the text at all — it simply identifies plausible word boundaries. Because of this limited scope, there is less need for application-specific custom parsers than there is for custom dictionaries. At present Postgres Pro provides just one built-in parser, which has been found to be useful for a wide range of applications.

The built-in parser is named `pg_catalog.default`. It recognizes 23 token types, shown in [Table 12.1](#).

Table 12.1. Default Parser's Token Types

Alias	Description	Example
<code>asciiword</code>	Word, all ASCII letters	elephant
<code>word</code>	Word, all letters	mañana

Alias	Description	Example
numword	Word, letters and digits	beta1
asciihword	Hyphenated word, all ASCII	up-to-date
hword	Hyphenated word, all letters	lógico-matemática
numhword	Hyphenated word, letters and digits	postgresql-beta1
hword_asciipart	Hyphenated word part, all ASCII	postgresql in the context post-gresql-beta1
hword_part	Hyphenated word part, all letters	lógico or matemática in the context lógico-matemática
hword_numpart	Hyphenated word part, letters and digits	beta1 in the context postgresql-beta1
email	Email address	foo@example.com
protocol	Protocol head	http://
url	URL	example.com/stuff/index.html
host	Host	example.com
url_path	URL path	/stuff/index.html, in the context of a URL
file	File or path name	/usr/local/foo.txt, if not within a URL
sfloat	Scientific notation	-1.234e56
float	Decimal notation	-1.234
int	Signed integer	-1234
uint	Unsigned integer	1234
version	Version number	8.3.0
tag	XML tag	
entity	XML entity	&
blank	Space symbols	(any whitespace or punctuation not otherwise recognized)

Note

The parser's notion of a “letter” is determined by the database's locale setting, specifically `lc_c-type`. Words containing only the basic ASCII letters are reported as a separate token type, since it is sometimes useful to distinguish them. In most European languages, token types `word` and `asciihword` should be treated alike.

`email` does not support all valid email characters as defined by [RFC 5322](#). Specifically, the only non-alphanumeric characters supported for email user names are period, dash, and underscore.

`tag` does not support all valid tag names as defined by [W3C Recommendation, XML](#). Specifically, the only tag names supported are those starting with an ASCII letter, underscore, or colon, and containing only letters, digits, hyphens, underscores, periods, and colons. `tag` also includes XML comments starting with `<!--` and ending with `-->`, and XML declarations (but note that this includes anything starting with `<?x` and ending with `>`).

It is possible for the parser to produce overlapping tokens from the same piece of text. As an example, a hyphenated word will be reported both as the entire word and as each component:

```
SELECT alias, description, token FROM ts_debug('foo-bar-beta1');
```

alias	description	token
numhword	Hyphenated word, letters and digits	foo-bar-beta1
hword_asciipart	Hyphenated word part, all ASCII	foo
blank	Space symbols	-
hword_asciipart	Hyphenated word part, all ASCII	bar
blank	Space symbols	-
hword_numpart	Hyphenated word part, letters and digits	beta1

This behavior is desirable since it allows searches to work for both the whole compound word and for components. Here is another instructive example:

```
SELECT alias, description, token FROM ts_debug('http://example.com/stuff/index.html');
```

alias	description	token
protocol	Protocol head	http://
url	URL	example.com/stuff/index.html
host	Host	example.com
url_path	URL path	/stuff/index.html

12.6. Dictionaries

Dictionaries are used to eliminate words that should not be considered in a search (*stop words*), and to *normalize* words so that different derived forms of the same word will match. A successfully normalized word is called a *lexeme*. Aside from improving search quality, normalization and removal of stop words reduce the size of the `tsvector` representation of a document, thereby improving performance. Normalization does not always have linguistic meaning and usually depends on application semantics.

Some examples of normalization:

- Linguistic — Ispell dictionaries try to reduce input words to a normalized form; stemmer dictionaries remove word endings
- URL locations can be canonicalized to make equivalent URLs match:
 - `http://www.pgsql.ru/db/mw/index.html`
 - `http://www.pgsql.ru/db/mw/`
 - `http://www.pgsql.ru/db/./db/mw/index.html`
- Color names can be replaced by their hexadecimal values, e.g., red, green, blue, magenta -> `FF0000`, `00FF00`, `0000FF`, `FF00FF`
- If indexing numbers, we can remove some fractional digits to reduce the range of possible numbers, so for example `3.14159265359`, `3.1415926`, `3.14` will be the same after normalization if only two digits are kept after the decimal point.

A dictionary is a program that accepts a token as input and returns:

- an array of lexemes if the input token is known to the dictionary (notice that one token can produce more than one lexeme)
- a single lexeme with the `TSL_FILTER` flag set, to replace the original token with a new token to be passed to subsequent dictionaries (a dictionary that does this is called a *filtering dictionary*)
- an empty array if the dictionary knows the token, but it is a stop word
- `NULL` if the dictionary does not recognize the input token

Postgres Pro provides predefined dictionaries for many languages. There are also several predefined templates that can be used to create new dictionaries with custom parameters. Each predefined dictionary template is described below. If no existing template is suitable, it is possible to create new ones; see the `contrib/` area of the Postgres Pro distribution for examples.

A text search configuration binds a parser together with a set of dictionaries to process the parser's output tokens. For each token type that the parser can return, a separate list of dictionaries is specified by the configuration. When a token of that type is found by the parser, each dictionary in the list is

consulted in turn, until some dictionary recognizes it as a known word. If it is identified as a stop word, or if no dictionary recognizes the token, it will be discarded and not indexed or searched for. Normally, the first dictionary that returns a non-NULL output determines the result, and any remaining dictionaries are not consulted; but a filtering dictionary can replace the given word with a modified word, which is then passed to subsequent dictionaries.

The general rule for configuring a list of dictionaries is to place first the most narrow, most specific dictionary, then the more general dictionaries, finishing with a very general dictionary, like a `Snowball` stemmer or `simple`, which recognizes everything. For example, for an astronomy-specific search (`astro_en` configuration) one could bind token type `asciiword` (ASCII word) to a synonym dictionary of astronomical terms, a general English dictionary and a `Snowball` English stemmer:

```
ALTER TEXT SEARCH CONFIGURATION astro_en
    ADD MAPPING FOR asciiword WITH astrosyn, english_ispell, english_stem;
```

A filtering dictionary can be placed anywhere in the list, except at the end where it'd be useless. Filtering dictionaries are useful to partially normalize words to simplify the task of later dictionaries. For example, a filtering dictionary could be used to remove accents from accented letters .

12.6.1. Stop Words

Stop words are words that are very common, appear in almost every document, and have no discrimination value. Therefore, they can be ignored in the context of full text searching. For example, every English text contains words like `a` and `the`, so it is useless to store them in an index. However, stop words do affect the positions in `tsvector`, which in turn affect ranking:

```
SELECT to_tsvector('english', 'in the list of stop words');
       to_tsvector
-----
'list':3 'stop':5 'word':6
```

The missing positions 1,2,4 are because of stop words. Ranks calculated for documents with and without stop words are quite different:

```
SELECT ts_rank_cd (to_tsvector('english', 'in the list of stop words'),
    to_tsquery('list & stop'));
       ts_rank_cd
-----
           0.05
```

```
SELECT ts_rank_cd (to_tsvector('english', 'list stop words'), to_tsquery('list &
    stop'));
       ts_rank_cd
-----
           0.1
```

It is up to the specific dictionary how it treats stop words. For example, `ispell` dictionaries first normalize words and then look at the list of stop words, while `Snowball` stemmers first check the list of stop words. The reason for the different behavior is an attempt to decrease noise.

12.6.2. Simple Dictionary

The `simple` dictionary template operates by converting the input token to lower case and checking it against a file of stop words. If it is found in the file then an empty array is returned, causing the token to be discarded. If not, the lower-cased form of the word is returned as the normalized lexeme. Alternatively, the dictionary can be configured to report non-stop-words as unrecognized, allowing them to be passed on to the next dictionary in the list.

Here is an example of a dictionary definition using the `simple` template:

```
CREATE TEXT SEARCH DICTIONARY public.simple_dict (
    TEMPLATE = pg_catalog.simple,
```

```
STOPWORDS = english
);
```

Here, `english` is the base name of a file of stop words. The file's full name will be `$$SHAREDIR/tsearch_data/english.stop`, where `$$SHAREDIR` means the Postgres Pro installation's shared-data directory, often `/usr/local/share/postgresql` (use `pg_config --sharedir` to determine it if you're not sure). The file format is simply a list of words, one per line. Blank lines and trailing spaces are ignored, and upper case is folded to lower case, but no other processing is done on the file contents.

Now we can test our dictionary:

```
SELECT ts_lexize('public.simple_dict', 'YeS');
ts_lexize
-----
{yes}
```

```
SELECT ts_lexize('public.simple_dict', 'The');
ts_lexize
-----
{}
```

We can also choose to return `NULL`, instead of the lower-cased word, if it is not found in the stop words file. This behavior is selected by setting the dictionary's `Accept` parameter to `false`. Continuing the example:

```
ALTER TEXT SEARCH DICTIONARY public.simple_dict ( Accept = false );
```

```
SELECT ts_lexize('public.simple_dict', 'YeS');
ts_lexize
-----
```

```
SELECT ts_lexize('public.simple_dict', 'The');
ts_lexize
-----
{}
```

With the default setting of `Accept = true`, it is only useful to place a `simple` dictionary at the end of a list of dictionaries, since it will never pass on any token to a following dictionary. Conversely, `Accept = false` is only useful when there is at least one following dictionary.

Caution

Most types of dictionaries rely on configuration files, such as files of stop words. These files *must* be stored in UTF-8 encoding. They will be translated to the actual database encoding, if that is different, when they are read into the server.

Caution

Normally, a database session will read a dictionary configuration file only once, when it is first used within the session. If you modify a configuration file and want to force existing sessions to pick up the new contents, issue an `ALTER TEXT SEARCH DICTIONARY` command on the dictionary. This can be a “dummy” update that doesn't actually change any parameter values.

12.6.3. Synonym Dictionary

This dictionary template is used to create dictionaries that replace a word with a synonym. Phrases are not supported (use the thesaurus template ([Section 12.6.4](#)) for that). A synonym dictionary can be used

to overcome linguistic problems, for example, to prevent an English stemmer dictionary from reducing the word “Paris” to “pari”. It is enough to have a `Paris pari` line in the synonym dictionary and put it before the `english_stem` dictionary. For example:

```
SELECT * FROM ts_debug('english', 'Paris');
  alias | description | token | dictionaries | dictionary | lexemes
-----+-----+-----+-----+-----+-----
asciiword | Word, all ASCII | Paris | {english_stem} | english_stem | {pari}

CREATE TEXT SEARCH DICTIONARY my_synonym (
    TEMPLATE = synonym,
    SYNONYMS = my_synonyms
);

ALTER TEXT SEARCH CONFIGURATION english
    ALTER MAPPING FOR asciiword
    WITH my_synonym, english_stem;

SELECT * FROM ts_debug('english', 'Paris');
  alias | description | token | dictionaries | dictionary | lexemes
-----+-----+-----+-----+-----+-----
+-----+
asciiword | Word, all ASCII | Paris | {my_synonym,english_stem} | my_synonym | {paris}
```

The only parameter required by the `synonym` template is `SYNONYMS`, which is the base name of its configuration file — `my_synonyms` in the above example. The file's full name will be `$SHAREDIR/tsearch_data/my_synonyms.syn` (where `$SHAREDIR` means the Postgres Pro installation's shared-data directory). The file format is just one line per word to be substituted, with the word followed by its synonym, separated by white space. Blank lines and trailing spaces are ignored.

The `synonym` template also has an optional parameter `CaseSensitive`, which defaults to `false`. When `CaseSensitive` is `false`, words in the synonym file are folded to lower case, as are input tokens. When it is `true`, words and tokens are not folded to lower case, but are compared as-is.

An asterisk (*) can be placed at the end of a synonym in the configuration file. This indicates that the synonym is a prefix. The asterisk is ignored when the entry is used in `to_tsvector()`, but when it is used in `to_tsquery()`, the result will be a query item with the prefix match marker (see [Section 12.3.2](#)). For example, suppose we have these entries in `$SHAREDIR/tsearch_data/synonym_sample.syn`:

```
postgres      pgsql
postgresql    pgsql
postgre pgsql
gogle  googl
indices index*
```

Then we will get these results:

```
mydb=# CREATE TEXT SEARCH DICTIONARY syn (template=synonym, synonyms='synonym_sample');
mydb=# SELECT ts_lexize('syn', 'indices');
 ts_lexize
-----
 {index}
(1 row)

mydb=# CREATE TEXT SEARCH CONFIGURATION tst (copy=simple);
mydb=# ALTER TEXT SEARCH CONFIGURATION tst ALTER MAPPING FOR asciiword WITH syn;
mydb=# SELECT to_tsvector('tst', 'indices');
 to_tsvector
-----
 'index':1
```

```
(1 row)
```

```
mydb=# SELECT to_tsquery('tst', 'indices');
 to_tsquery
```

```
-----
 'index':*
(1 row)
```

```
mydb=# SELECT 'indexes are very useful'::tsvector;
 tsvector
```

```
-----
 'are' 'indexes' 'useful' 'very'
(1 row)
```

```
mydb=# SELECT 'indexes are very useful'::tsvector @@ to_tsquery('tst', 'indices');
 ?column?
```

```
-----
 t
(1 row)
```

12.6.4. Thesaurus Dictionary

A thesaurus dictionary (sometimes abbreviated as TZ) is a collection of words that includes information about the relationships of words and phrases, i.e., broader terms (BT), narrower terms (NT), preferred terms, non-preferred terms, related terms, etc.

Basically a thesaurus dictionary replaces all non-preferred terms by one preferred term and, optionally, preserves the original terms for indexing as well. Postgres Pro's current implementation of the thesaurus dictionary is an extension of the synonym dictionary with added *phrase* support. A thesaurus dictionary requires a configuration file of the following format:

```
# this is a comment
sample word(s) : indexed word(s)
more sample word(s) : more indexed word(s)
...
```

where the colon (:) symbol acts as a delimiter between a phrase and its replacement.

A thesaurus dictionary uses a *subdictionary* (which is specified in the dictionary's configuration) to normalize the input text before checking for phrase matches. It is only possible to select one subsdictionary. An error is reported if the subsdictionary fails to recognize a word. In that case, you should remove the use of the word or teach the subsdictionary about it. You can place an asterisk (*) at the beginning of an indexed word to skip applying the subsdictionary to it, but all sample words *must* be known to the subsdictionary.

The thesaurus dictionary chooses the longest match if there are multiple phrases matching the input, and ties are broken by using the last definition.

Specific stop words recognized by the subsdictionary cannot be specified; instead use ? to mark the location where any stop word can appear. For example, assuming that `a` and `the` are stop words according to the subsdictionary:

```
? one ? two : swsw
```

matches `a one the two and the one a two`; both would be replaced by `swsw`.

Since a thesaurus dictionary has the capability to recognize phrases it must remember its state and interact with the parser. A thesaurus dictionary uses these assignments to check if it should handle the next word or stop accumulation. The thesaurus dictionary must be configured carefully. For example, if the thesaurus dictionary is assigned to handle only the `asciiword` token, then a thesaurus dictionary definition like `one 7` will not work since token type `uint` is not assigned to the thesaurus dictionary.

Caution

Thesauruses are used during indexing so any change in the thesaurus dictionary's parameters *requires* reindexing. For most other dictionary types, small changes such as adding or removing stopwords does not force reindexing.

12.6.4.1. Thesaurus Configuration

To define a new thesaurus dictionary, use the `thesaurus` template. For example:

```
CREATE TEXT SEARCH DICTIONARY thesaurus_simple (  
    TEMPLATE = thesaurus,  
    DictFile = mythesaurus,  
    Dictionary = pg_catalog.english_stem  
);
```

Here:

- `thesaurus_simple` is the new dictionary's name
- `mythesaurus` is the base name of the thesaurus configuration file. (Its full name will be `$SHAREDIR/tsearch_data/mythesaurus.ths`, where `$SHAREDIR` means the installation shared-data directory.)
- `pg_catalog.english_stem` is the subdictionary (here, a Snowball English stemmer) to use for thesaurus normalization. Notice that the subdictionary will have its own configuration (for example, stop words), which is not shown here.

Now it is possible to bind the thesaurus dictionary `thesaurus_simple` to the desired token types in a configuration, for example:

```
ALTER TEXT SEARCH CONFIGURATION russian  
    ALTER MAPPING FOR asciiword, asciihword, hword_asciipart  
    WITH thesaurus_simple;
```

12.6.4.2. Thesaurus Example

Consider a simple astronomical thesaurus `thesaurus_astro`, which contains some astronomical word combinations:

```
supernovae stars : sn  
crab nebulae : crab
```

Below we create a dictionary and bind some token types to an astronomical thesaurus and English stemmer:

```
CREATE TEXT SEARCH DICTIONARY thesaurus_astro (  
    TEMPLATE = thesaurus,  
    DictFile = thesaurus_astro,  
    Dictionary = english_stem  
);  
  
ALTER TEXT SEARCH CONFIGURATION russian  
    ALTER MAPPING FOR asciiword, asciihword, hword_asciipart  
    WITH thesaurus_astro, english_stem;
```

Now we can see how it works. `ts_lexize` is not very useful for testing a thesaurus, because it treats its input as a single token. Instead we can use `plainto_tsquery` and `to_tsvector` which will break their input strings into multiple tokens:

```
SELECT plainto_tsquery('supernova star');  
plainto_tsquery  
-----  
'sn'
```

```
SELECT to_tsvector('supernova star');
to_tsvector
-----
'sn':1
```

In principle, one can use `to_tsquery` if you quote the argument:

```
SELECT to_tsquery('supernova star');
to_tsquery
-----
'sn'
```

Notice that `supernova star` matches `supernovae stars` in `thesaurus_astro` because we specified the `english_stem` stemmer in the thesaurus definition. The stemmer removed the `e` and `s`.

To index the original phrase as well as the substitute, just include it in the right-hand part of the definition:

```
supernovae stars : sn supernovae stars

SELECT plainto_tsquery('supernova star');
plainto_tsquery
-----
'sn' & 'supernova' & 'star'
```

12.6.5. Ispell Dictionary

The Ispell dictionary template supports *morphological dictionaries*, which can normalize many different linguistic forms of a word into the same lexeme. For example, an English Ispell dictionary can match all declensions and conjugations of the search term `bank`, e.g., `banking`, `banked`, `banks`, `banks'`, and `bank's`.

The standard Postgres Pro distribution does not include any Ispell configuration files. Dictionaries for a large number of languages are available from [Ispell](#). Also, some more modern dictionary file formats are supported — [MySpell](#) (OO < 2.0.1) and [Hunspell](#) (OO >= 2.0.2). A large list of dictionaries is available on the [OpenOffice Wiki](#).

To create an Ispell dictionary perform these steps:

- download dictionary configuration files. OpenOffice extension files have the `.oxt` extension. It is necessary to extract `.aff` and `.dic` files, change extensions to `.affix` and `.dict`. For some dictionary files it is also needed to convert characters to the UTF-8 encoding with commands (for example, for a Norwegian language dictionary):

```
iconv -f ISO_8859-1 -t UTF-8 -o nn_no.affix nn_NO.aff
iconv -f ISO_8859-1 -t UTF-8 -o nn_no.dict nn_NO.dic
```

- copy files to the `$$SHAREDIR/tsearch_data` directory
- load files into Postgres Pro with the following command:

```
CREATE TEXT SEARCH DICTIONARY english_hunspell (
    TEMPLATE = ispell,
    DictFile = en_us,
    AffFile = en_us,
    Stopwords = english);
```

Here, `DictFile`, `AffFile`, and `StopWords` specify the base names of the dictionary, affixes, and stop-words files. The stop-words file has the same format explained above for the `simple` dictionary type. The format of the other files is not specified here but is available from the above-mentioned web sites.

Ispell dictionaries usually recognize a limited set of words, so they should be followed by another broader dictionary; for example, a Snowball dictionary, which recognizes everything.

The `.affix` file of Ispell has the following structure:

```
prefixes
flag *A:
.          > RE          # As in enter > reenter
suffixes
flag T:
E          > ST          # As in late > latest
[^AEIOU]Y  > -Y, IEST    # As in dirty > dirtiest
[AEIOU]Y   > EST        # As in gray > grayest
[^EY]      > EST        # As in small > smallest
```

And the `.dict` file has the following structure:

```
lapse/ADGRS
lard/DGRS
large/PRTY
lark/MRS
```

Format of the `.dict` file is:

```
basic_form/affix_class_name
```

In the `.affix` file every affix flag is described in the following format:

```
condition > [-stripping_letters,] adding_affix
```

Here, condition has a format similar to the format of regular expressions. It can use groupings [...] and [^...]. For example, [AEIOU]Y means that the last letter of the word is "y" and the penultimate letter is "a", "e", "i", "o" or "u". [^EY] means that the last letter is neither "e" nor "y".

Ispell dictionaries support splitting compound words; a useful feature. Notice that the affix file should specify a special flag using the `compoundwords controlled` statement that marks dictionary words that can participate in compound formation:

```
compoundwords controlled z
```

Here are some examples for the Norwegian language:

```
SELECT ts_lexize('norwegian_ispell', 'overbuljongtarningpakkmasterassistent');
{over,buljong,terning,pakk,mester,assistent}
SELECT ts_lexize('norwegian_ispell', 'sjokoladefabrikk');
{sjokoladefabrikk,sjokolade,fabrikk}
```

MySpell format is a subset of Hunspell. The `.affix` file of Hunspell has the following structure:

```
PFX A Y 1
PFX A 0 re .
SFX T N 4
SFX T 0 st e
SFX T y iest [^aeiou]y
SFX T 0 est [aeiou]y
SFX T 0 est [^ey]
```

The first line of an affix class is the header. Fields of an affix rules are listed after the header:

- parameter name (PFX or SFX)
- flag (name of the affix class)
- stripping characters from beginning (at prefix) or end (at suffix) of the word
- adding affix
- condition that has a format similar to the format of regular expressions.

The `.dict` file looks like the `.dict` file of Ispell:

```
larder/M
```

```
lardy/RT
large/RSPMYT
largehearted
```

Note

MySpell does not support compound words. Hunspell has sophisticated support for compound words. At present, Postgres Pro implements only the basic compound word operations of Hunspell.

12.6.6. Snowball Dictionary

The Snowball dictionary template is based on a project by Martin Porter, inventor of the popular Porter's stemming algorithm for the English language. Snowball now provides stemming algorithms for many languages (see the [Snowball site](#) for more information). Each algorithm understands how to reduce common variant forms of words to a base, or stem, spelling within its language. A Snowball dictionary requires a `language` parameter to identify which stemmer to use, and optionally can specify a `stopword` file name that gives a list of words to eliminate. (Postgres Pro's standard stopwords lists are also provided by the Snowball project.) For example, there is a built-in definition equivalent to

```
CREATE TEXT SEARCH DICTIONARY english_stem (
    TEMPLATE = snowball,
    Language = english,
    StopWords = english
);
```

The stopwords file format is the same as already explained.

A Snowball dictionary recognizes everything, whether or not it is able to simplify the word, so it should be placed at the end of the dictionary list. It is useless to have it before any other dictionary because a token will never pass through it to the next dictionary.

12.7. Configuration Example

A text search configuration specifies all options necessary to transform a document into a `tsvector`: the parser to use to break text into tokens, and the dictionaries to use to transform each token into a lexeme. Every call of `to_tsvector` or `to_tsquery` needs a text search configuration to perform its processing. The configuration parameter `default_text_search_config` specifies the name of the default configuration, which is the one used by text search functions if an explicit configuration parameter is omitted. It can be set in `postgresql.conf`, or set for an individual session using the `SET` command.

Several predefined text search configurations are available, and you can create custom configurations easily. To facilitate management of text search objects, a set of SQL commands is available, and there are several `psql` commands that display information about text search objects ([Section 12.10](#)).

As an example we will create a configuration `pg`, starting by duplicating the built-in `english` configuration:

```
CREATE TEXT SEARCH CONFIGURATION public.pg ( COPY = pg_catalog.english );
```

We will use a Postgres Pro-specific synonym list and store it in `$SHAREDIR/tsearch_data/pg_dict.syn`. The file contents look like:

```
postgres    pg
pgsql       pg
postgresql  pg
```

We define the synonym dictionary like this:

```
CREATE TEXT SEARCH DICTIONARY pg_dict (
    TEMPLATE = synonym,
```

```
SYNONYMS = pg_dict
);
```

Next we register the Ispell dictionary `english_ispell`, which has its own configuration files:

```
CREATE TEXT SEARCH DICTIONARY english_ispell (
    TEMPLATE = ispell,
    DictFile = english,
    AffFile = english,
    StopWords = english
);
```

Now we can set up the mappings for words in configuration `pg`:

```
ALTER TEXT SEARCH CONFIGURATION pg
    ALTER MAPPING FOR asciiword, asciihword, hword_asciipart,
                        word, hword, hword_part
    WITH pg_dict, english_ispell, english_stem;
```

We choose not to index or search some token types that the built-in configuration does handle:

```
ALTER TEXT SEARCH CONFIGURATION pg
    DROP MAPPING FOR email, url, url_path, sfloat, float;
```

Now we can test our configuration:

```
SELECT * FROM ts_debug('public.pg', '
Postgres Pro, the highly scalable, SQL compliant, open source object-relational
database management system, is now undergoing beta testing of the next
version of our software.
');
```

The next step is to set the session to use the new configuration, which was created in the `public` schema:

```
=> \dF
      List of text search configurations
 Schema | Name | Description
-----+-----+-----
 public | pg   |

SET default_text_search_config = 'public.pg';
SET

SHOW default_text_search_config;
 default_text_search_config
-----
 public.pg
```

12.8. Testing and Debugging Text Search

The behavior of a custom text search configuration can easily become confusing. The functions described in this section are useful for testing text search objects. You can test a complete configuration, or test parsers and dictionaries separately.

12.8.1. Configuration Testing

The function `ts_debug` allows easy testing of a text search configuration.

```
ts_debug([ config regconfig, ] document text,
         OUT alias text,
         OUT description text,
         OUT token text,
```

```

OUT dictionaries regdictionary[],
OUT dictionary regdictionary,
OUT lexemes text[])
returns setof record

```

`ts_debug` displays information about every token of *document* as produced by the parser and processed by the configured dictionaries. It uses the configuration specified by *config*, or *default_text_search_config* if that argument is omitted.

`ts_debug` returns one row for each token identified in the text by the parser. The columns returned are

- *alias* text — short name of the token type
- *description* text — description of the token type
- *token* text — text of the token
- *dictionaries* regdictionary[] — the dictionaries selected by the configuration for this token type
- *dictionary* regdictionary — the dictionary that recognized the token, or NULL if none did
- *lexemes* text[] — the lexeme(s) produced by the dictionary that recognized the token, or NULL if none did; an empty array ({}) means it was recognized as a stop word

Here is a simple example:

```

SELECT * FROM ts_debug('english', 'a fat  cat sat on a mat - it ate a fat rats');

```

alias	description	token	dictionaries	dictionary	lexemes
asciiword	Word, all ASCII	a	{english_stem}	english_stem	{}
blank	Space symbols		{}		
asciiword	Word, all ASCII	fat	{english_stem}	english_stem	{fat}
blank	Space symbols		{}		
asciiword	Word, all ASCII	cat	{english_stem}	english_stem	{cat}
blank	Space symbols		{}		
asciiword	Word, all ASCII	sat	{english_stem}	english_stem	{sat}
blank	Space symbols		{}		
asciiword	Word, all ASCII	on	{english_stem}	english_stem	{}
blank	Space symbols		{}		
asciiword	Word, all ASCII	a	{english_stem}	english_stem	{}
blank	Space symbols		{}		
asciiword	Word, all ASCII	mat	{english_stem}	english_stem	{mat}
blank	Space symbols		{}		
blank	Space symbols	-	{}		
asciiword	Word, all ASCII	it	{english_stem}	english_stem	{}
blank	Space symbols		{}		
asciiword	Word, all ASCII	ate	{english_stem}	english_stem	{ate}
blank	Space symbols		{}		
asciiword	Word, all ASCII	a	{english_stem}	english_stem	{}
blank	Space symbols		{}		
asciiword	Word, all ASCII	fat	{english_stem}	english_stem	{fat}
blank	Space symbols		{}		
asciiword	Word, all ASCII	rats	{english_stem}	english_stem	{rat}

For a more extensive demonstration, we first create a `public.english` configuration and Ispell dictionary for the English language:

```

CREATE TEXT SEARCH CONFIGURATION public.english ( COPY = pg_catalog.english );

CREATE TEXT SEARCH DICTIONARY english_ispell (
    TEMPLATE = ispell,
    DictFile = english,
    AffFile = english,
    StopWords = english
);

```

```

ALTER TEXT SEARCH CONFIGURATION public.english
  ALTER MAPPING FOR asciiword WITH english_ispell, english_stem;

SELECT * FROM ts_debug('public.english', 'The Brightest supernovaes');

```

alias	description	token	dictionaries
dictionary	lexemes		
asciiword	Word, all ASCII	The	{english_ispell,english_stem}
english_ispell	{}		
blank	Space symbols		{}
asciiword	Word, all ASCII	Brightest	{english_ispell,english_stem}
english_ispell	{bright}		
blank	Space symbols		{}
asciiword	Word, all ASCII	supernovaes	{english_ispell,english_stem}
english_stem	{supernova}		

In this example, the word `Brightest` was recognized by the parser as an ASCII word (alias `asciiword`). For this token type the dictionary list is `english_ispell` and `english_stem`. The word was recognized by `english_ispell`, which reduced it to the noun `bright`. The word `supernovaes` is unknown to the `english_ispell` dictionary so it was passed to the next dictionary, and, fortunately, was recognized (in fact, `english_stem` is a Snowball dictionary which recognizes everything; that is why it was placed at the end of the dictionary list).

The word `The` was recognized by the `english_ispell` dictionary as a stop word ([Section 12.6.1](#)) and will not be indexed. The spaces are discarded too, since the configuration provides no dictionaries at all for them.

You can reduce the width of the output by explicitly specifying which columns you want to see:

```

SELECT alias, token, dictionary, lexemes
FROM ts_debug('public.english', 'The Brightest supernovaes');

```

alias	token	dictionary	lexemes
asciiword	The	english_ispell	{}
blank			
asciiword	Brightest	english_ispell	{bright}
blank			
asciiword	supernovaes	english_stem	{supernova}

12.8.2. Parser Testing

The following functions allow direct testing of a text search parser.

```

ts_parse(parser_name text, document text,
         OUT tokid integer, OUT token text) returns setof record
ts_parse(parser_oid oid, document text,
         OUT tokid integer, OUT token text) returns setof record

```

`ts_parse` parses the given *document* and returns a series of records, one for each token produced by parsing. Each record includes a `tokid` showing the assigned token type and a `token` which is the text of the token. For example:

```

SELECT * FROM ts_parse('default', '123 - a number');

```

tokid	token
22	123
12	-

```
12 | -
  1 | a
12 |
  1 | number
```

```
ts_token_type(parser_name text, OUT tokid integer,
              OUT alias text, OUT description text) returns setof record
ts_token_type(parser_oid oid, OUT tokid integer,
              OUT alias text, OUT description text) returns setof record
```

`ts_token_type` returns a table which describes each type of token the specified parser can recognize. For each token type, the table gives the integer `tokid` that the parser uses to label a token of that type, the alias that names the token type in configuration commands, and a short description. For example:

```
SELECT * FROM ts_token_type('default');
 tokid |      alias      | description
-----+-----+-----
    1 | asciiword       | Word, all ASCII
    2 | word            | Word, all letters
    3 | numword         | Word, letters and digits
    4 | email           | Email address
    5 | url             | URL
    6 | host            | Host
    7 | sfloat          | Scientific notation
    8 | version         | Version number
    9 | hword_numpart   | Hyphenated word part, letters and digits
   10 | hword_part      | Hyphenated word part, all letters
   11 | hword_asciipart | Hyphenated word part, all ASCII
   12 | blank           | Space symbols
   13 | tag             | XML tag
   14 | protocol        | Protocol head
   15 | numhword        | Hyphenated word, letters and digits
   16 | asciihword      | Hyphenated word, all ASCII
   17 | hword           | Hyphenated word, all letters
   18 | url_path        | URL path
   19 | file            | File or path name
   20 | float           | Decimal notation
   21 | int             | Signed integer
   22 | uint            | Unsigned integer
   23 | entity          | XML entity
```

12.8.3. Dictionary Testing

The `ts_lexize` function facilitates dictionary testing.

```
ts_lexize(dict regdictionary, token text) returns text[]
```

`ts_lexize` returns an array of lexemes if the input `token` is known to the dictionary, or an empty array if the token is known to the dictionary but it is a stop word, or `NULL` if it is an unknown word.

Examples:

```
SELECT ts_lexize('english_stem', 'stars');
 ts_lexize
-----
 {star}
```

```
SELECT ts_lexize('english_stem', 'a');
 ts_lexize
-----
```


{}

Note

The `ts_lexize` function expects a single *token*, not text. Here is a case where this can be confusing:

```
SELECT ts_lexize('thesaurus_astro', 'supernovae stars') is null;
?column?
```

```
-----
t
```

The thesaurus dictionary `thesaurus_astro` does know the phrase `supernovae stars`, but `ts_lexize` fails since it does not parse the input text but treats it as a single token. Use `plainto_tsquery` or `to_tsvector` to test thesaurus dictionaries, for example:

```
SELECT plainto_tsquery('supernovae stars');
plainto_tsquery
-----
'sn'
```

12.9. Preferred Index Types for Text Search

The following kinds of indexes can be used to speed up full text searches: [GIN](#) and [GiST](#). Note that indexes are not mandatory for full text searching, but in cases where a column is searched on a regular basis, an index is usually desirable.

To create such an index, do one of:

```
CREATE INDEX name ON table USING GIN (column);
```

Creates a GIN (Generalized Inverted Index)-based index. The *column* must be of `tsvector` type.

```
CREATE INDEX name ON table USING RUM (column);
```

Creates a RUM-based index. The *column* must be of `tsvector` type.

```
CREATE INDEX name ON table USING GIST (column [ { DEFAULT | tsvector_ops } (siglen =
number) ] );
```

Creates a GiST (Generalized Search Tree)-based index. The *column* can be of `tsvector` or `tsquery` type. Optional integer parameter `siglen` determines signature length in bytes (see below for details).

GIN indexes and RUM indexes, based on the GIN access method code, are the preferred text search index type. As inverted indexes, they contain an index entry for each word (lexeme), with a compressed list of matching locations. Multi-word searches can find the first match, then use the index to remove rows that are lacking additional words. GIN indexes store only the words (lexemes) of `tsvector` values, and not their weight labels. Thus a table row recheck is needed when using a query that involves weights.

A GiST index is *lossy*, meaning that the index might produce false matches, and it is necessary to check the actual table row to eliminate such false matches. (Postgres Pro does this automatically when needed.) GiST indexes are lossy because each document is represented in the index by a fixed-length signature. The signature length in bytes is determined by the value of the optional integer parameter `siglen`. The default signature length (when `siglen` is not specified) is 124 bytes, the maximum signature length is 2024 bytes. The signature is generated by hashing each word into a single bit in an *n*-bit string, with all these bits OR-ed together to produce an *n*-bit document signature. When two words hash to the same bit position there will be a false match. If all words in the query have matches (real or false) then the table row must be retrieved to see if the match is correct. Longer signatures lead to a more precise search (scanning a smaller fraction of the index and fewer heap pages), at the cost of a larger index.

A GiST index can be covering, i.e., use the `INCLUDE` clause. Included columns can have data types without any GiST operator class. Included attributes will be stored uncompressed.

Lossiness causes performance degradation due to unnecessary fetches of table records that turn out to be false matches. Since random access to table records is slow, this limits the usefulness of GiST indexes. The likelihood of false matches depends on several factors, in particular the number of unique words, so using dictionaries to reduce this number is recommended.

Note that GIN index build time can often be improved by increasing `maintenance_work_mem`, while GiST index build time is not sensitive to that parameter.

Partitioning of big collections and the proper use of GIN, RUM, and GiST indexes allows the implementation of very fast searches with online update. Partitioning can be done at the database level using table inheritance, or by distributing documents over servers and collecting external search results, e.g., via [Foreign Data](#) access. The latter is possible because ranking functions use only local information.

12.10. psql Support

Information about text search configuration objects can be obtained in psql using a set of commands:

```
\dF{d,p,t}[+] [PATTERN]
```

An optional `+` produces more details.

The optional parameter `PATTERN` can be the name of a text search object, optionally schema-qualified. If `PATTERN` is omitted then information about all visible objects will be displayed. `PATTERN` can be a regular expression and can provide *separate* patterns for the schema and object names. The following examples illustrate this:

```
=> \dF *fulltext*
      List of text search configurations
 Schema | Name          | Description
-----+-----+-----
 public | fulltext_cfg |

=> \dF *.fulltext*
      List of text search configurations
 Schema | Name          | Description
-----+-----+-----
 fulltext | fulltext_cfg |
 public   | fulltext_cfg |
```

The available commands are:

```
\dF[+] [PATTERN]
```

List text search configurations (add `+` for more detail).

```
=> \dF russian
      List of text search configurations
 Schema | Name          | Description
-----+-----+-----
 pg_catalog | russian | configuration for russian language

=> \dF+ russian
Text search configuration "pg_catalog.russian"
Parser: "pg_catalog.default"
 Token | Dictionaries
-----+-----
 asciihword | english_stem
 asciiword  | english_stem
 email      | simple
 file       | simple
 float      | simple
 host       | simple
```

hword		russian_stem
hword_asciipart		english_stem
hword_numpart		simple
hword_part		russian_stem
int		simple
numhword		simple
numword		simple
sfloat		simple
uint		simple
url		simple
url_path		simple
version		simple
word		russian_stem

\dFd[+] [PATTERN]

List text search dictionaries (add + for more detail).

=> \dFd

Schema	Name	Description
pg_catalog	arabic_stem	snowball stemmer for arabic language
pg_catalog	armenian_stem	snowball stemmer for armenian language
pg_catalog	basque_stem	snowball stemmer for basque language
pg_catalog	catalan_stem	snowball stemmer for catalan language
pg_catalog	danish_stem	snowball stemmer for danish language
pg_catalog	dutch_stem	snowball stemmer for dutch language
pg_catalog	english_stem	snowball stemmer for english language
pg_catalog	finnish_stem	snowball stemmer for finnish language
pg_catalog	french_stem	snowball stemmer for french language
pg_catalog	german_stem	snowball stemmer for german language
pg_catalog	greek_stem	snowball stemmer for greek language
pg_catalog	hindi_stem	snowball stemmer for hindi language
pg_catalog	hungarian_stem	snowball stemmer for hungarian language
pg_catalog	indonesian_stem	snowball stemmer for indonesian language
pg_catalog	irish_stem	snowball stemmer for irish language
pg_catalog	italian_stem	snowball stemmer for italian language
pg_catalog	lithuanian_stem	snowball stemmer for lithuanian language
pg_catalog	nepali_stem	snowball stemmer for nepali language
pg_catalog	norwegian_stem	snowball stemmer for norwegian language
pg_catalog	portuguese_stem	snowball stemmer for portuguese language
pg_catalog	romanian_stem	snowball stemmer for romanian language
pg_catalog	russian_stem	snowball stemmer for russian language
pg_catalog	serbian_stem	snowball stemmer for serbian language
pg_catalog	simple	simple dictionary: just lower case and check for stopword
pg_catalog	spanish_stem	snowball stemmer for spanish language
pg_catalog	swedish_stem	snowball stemmer for swedish language
pg_catalog	tamil_stem	snowball stemmer for tamil language
pg_catalog	turkish_stem	snowball stemmer for turkish language
pg_catalog	yiddish_stem	snowball stemmer for yiddish language

\dFp[+] [PATTERN]

List text search parsers (add + for more detail).

=> \dFp

List of text search parsers

Schema	Name	Description
pg_catalog	default	default word parser
=> \dFp+		
Text search parser "pg_catalog.default"		
Method	Function	Description
Start parse	prsd_start	
Get next token	prsd_nexttoken	
End parse	prsd_end	
Get headline	prsd_headline	
Get token types	prsd_lextype	

Token types for parser "pg_catalog.default"	
Token name	Description
asciihword	Hyphenated word, all ASCII
asciiword	Word, all ASCII
blank	Space symbols
email	Email address
entity	XML entity
file	File or path name
float	Decimal notation
host	Host
hword	Hyphenated word, all letters
hword_asciipart	Hyphenated word part, all ASCII
hword_numpart	Hyphenated word part, letters and digits
hword_part	Hyphenated word part, all letters
int	Signed integer
numhword	Hyphenated word, letters and digits
numword	Word, letters and digits
protocol	Protocol head
sfloat	Scientific notation
tag	XML tag
uint	Unsigned integer
url	URL
url_path	URL path
version	Version number
word	Word, all letters
(23 rows)	

\dFt[+] [PATTERN]

List text search templates (add + for more detail).

=> \dFt

List of text search templates		
Schema	Name	Description
pg_catalog	ispell	ispell dictionary
pg_catalog	simple	simple dictionary: just lower case and check for stopword
pg_catalog	snowball	snowball stemmer
pg_catalog	synonym	synonym dictionary: replace word by its synonym
pg_catalog	thesaurus	thesaurus dictionary: phrase by phrase substitution

12.11. Limitations

The current limitations of Postgres Pro's text search features are:

- The length of each lexeme must be less than 2 kilobytes

- The length of a `tsvector` (lexemes + positions) must be less than 1 megabyte
- The number of lexemes must be less than 2^{64}
- Position values in `tsvector` must be greater than 0 and no more than 16,383
- The match distance in a `<N>` (FOLLOWED BY) `tsquery` operator cannot be more than 16,384
- No more than 256 positions per lexeme
- The number of nodes (lexemes + operators) in a `tsquery` must be less than 32,768

For comparison, the PostgreSQL 8.1 documentation contained 10,441 unique words, a total of 335,420 words, and the most frequent word “postgresql” was mentioned 6,127 times in 655 documents.

Another example — the PostgreSQL mailing list archives contained 910,989 unique words with 57,491,343 lexemes in 461,020 messages.

Chapter 13. Concurrency Control

This chapter describes the behavior of the Postgres Pro database system when two or more sessions try to access the same data at the same time. The goals in that situation are to allow efficient access for all sessions while maintaining strict data integrity. Every developer of database applications should be familiar with the topics covered in this chapter.

13.1. Introduction

Postgres Pro provides a rich set of tools for developers to manage concurrent access to data. Internally, data consistency is maintained by using a multiversion model (Multiversion Concurrency Control, MVCC). This means that each SQL statement sees a snapshot of data (a *database version*) as it was some time ago, regardless of the current state of the underlying data. This prevents statements from viewing inconsistent data produced by concurrent transactions performing updates on the same data rows, providing *transaction isolation* for each database session. MVCC, by eschewing the locking methodologies of traditional database systems, minimizes lock contention in order to allow for reasonable performance in multiuser environments.

The main advantage of using the MVCC model of concurrency control rather than locking is that in MVCC locks acquired for querying (reading) data do not conflict with locks acquired for writing data, and so reading never blocks writing and writing never blocks reading. Postgres Pro maintains this guarantee even when providing the strictest level of transaction isolation through the use of an innovative *Serializable Snapshot Isolation* (SSI) level.

Table- and row-level locking facilities are also available in Postgres Pro for applications which don't generally need full transaction isolation and prefer to explicitly manage particular points of conflict. However, proper use of MVCC will generally provide better performance than locks. In addition, application-defined advisory locks provide a mechanism for acquiring locks that are not tied to a single transaction.

13.2. Transaction Isolation

The SQL standard defines four levels of transaction isolation. The most strict is Serializable, which is defined by the standard in a paragraph which says that any concurrent execution of a set of Serializable transactions is guaranteed to produce the same effect as running them one at a time in some order. The other three levels are defined in terms of phenomena, resulting from interaction between concurrent transactions, which must not occur at each level. The standard notes that due to the definition of Serializable, none of these phenomena are possible at that level. (This is hardly surprising -- if the effect of the transactions must be consistent with having been run one at a time, how could you see any phenomena caused by interactions?)

The phenomena which are prohibited at various levels are:

dirty read

A transaction reads data written by a concurrent uncommitted transaction.

nonrepeatable read

A transaction re-reads data it has previously read and finds that data has been modified by another transaction (that committed since the initial read).

phantom read

A transaction re-executes a query returning a set of rows that satisfy a search condition and finds that the set of rows satisfying the condition has changed due to another recently-committed transaction.

serialization anomaly

The result of successfully committing a group of transactions is inconsistent with all possible orderings of running those transactions one at a time.

The SQL standard and Postgres Pro-implemented transaction isolation levels are described in [Table 13.1](#).

Table 13.1. Transaction Isolation Levels

Isolation Level	Dirty Read	Nonrepeatable Read	Phantom Read	Serialization Anomaly
Read uncommitted	Allowed, but not in PG	Possible	Possible	Possible
Read committed	Not possible	Possible	Possible	Possible
Repeatable read	Not possible	Not possible	Allowed, but not in PG	Possible
Serializable	Not possible	Not possible	Not possible	Not possible

In Postgres Pro, you can request any of the four standard transaction isolation levels, but internally only three distinct isolation levels are implemented, i.e., Postgres Pro's Read Uncommitted mode behaves like Read Committed. This is because it is the only sensible way to map the standard isolation levels to Postgres Pro's multiversion concurrency control architecture.

The table also shows that Postgres Pro's Repeatable Read implementation does not allow phantom reads. This is acceptable under the SQL standard because the standard specifies which anomalies must *not* occur at certain isolation levels; higher guarantees are acceptable. The behavior of the available isolation levels is detailed in the following subsections.

To set the transaction isolation level of a transaction, use the command [SET TRANSACTION](#).

Important

Some Postgres Pro data types and functions have special rules regarding transactional behavior. In particular, changes made to a sequence (and therefore the counter of a column declared using `serial`) are immediately visible to all other transactions and are not rolled back if the transaction that made the changes aborts. See [Section 9.1.7](#) and [Section 8.1.4](#).

13.2.1. Read Committed Isolation Level

Read Committed is the default isolation level in Postgres Pro. When a transaction uses this isolation level, a `SELECT` query (without a `FOR UPDATE/SHARE` clause) sees only data committed before the query began; it never sees either uncommitted data or changes committed by concurrent transactions during the query's execution. In effect, a `SELECT` query sees a snapshot of the database as of the instant the query begins to run. However, `SELECT` does see the effects of previous updates executed within its own transaction, even though they are not yet committed. Also note that two successive `SELECT` commands can see different data, even though they are within a single transaction, if other transactions commit changes after the first `SELECT` starts and before the second `SELECT` starts.

`UPDATE`, `DELETE`, `SELECT FOR UPDATE`, and `SELECT FOR SHARE` commands behave the same as `SELECT` in terms of searching for target rows: they will only find target rows that were committed as of the command start time. However, such a target row might have already been updated (or deleted or locked) by another concurrent transaction by the time it is found. In this case, the would-be updater will wait for the first updating transaction to commit or roll back (if it is still in progress). If the first updater rolls back, then its effects are negated and the second updater can proceed with updating the originally found row. If the first updater commits, the second updater will ignore the row if the first updater deleted it, otherwise it will attempt to apply its operation to the updated version of the row. The search condition of the command (the `WHERE` clause) is re-evaluated to see if the updated version of the row still matches the search condition. If so, the second updater proceeds with its operation using the updated version of the row. In the case of `SELECT FOR UPDATE` and `SELECT FOR SHARE`, this means it is the updated version of the row that is locked and returned to the client.

`INSERT` with an `ON CONFLICT DO UPDATE` clause behaves similarly. In Read Committed mode, each row proposed for insertion will either insert or update. Unless there are unrelated errors, one of those two outcomes is guaranteed. If a conflict originates in another transaction whose effects are not yet visible to the `INSERT`, the `UPDATE` clause will affect that row, even though possibly *no* version of that row is conventionally visible to the command.

`INSERT` with an `ON CONFLICT DO NOTHING` clause may have insertion not proceed for a row due to the outcome of another transaction whose effects are not visible to the `INSERT` snapshot. Again, this is only the case in Read Committed mode.

`MERGE` allows the user to specify various combinations of `INSERT`, `UPDATE` and `DELETE` subcommands. A `MERGE` command with both `INSERT` and `UPDATE` subcommands looks similar to `INSERT` with an `ON CONFLICT DO UPDATE` clause but does not guarantee that either `INSERT` or `UPDATE` will occur. If `MERGE` attempts an `UPDATE` or `DELETE` and the row is concurrently updated but the join condition still passes for the current target and the current source tuple, then `MERGE` will behave the same as the `UPDATE` or `DELETE` commands and perform its action on the updated version of the row. However, because `MERGE` can specify several actions and they can be conditional, the conditions for each action are re-evaluated on the updated version of the row, starting from the first action, even if the action that had originally matched appears later in the list of actions. On the other hand, if the row is concurrently updated so that the join condition fails, then `MERGE` will evaluate the command's `NOT MATCHED BY SOURCE` and `NOT MATCHED [BY TARGET]` actions next, and execute the first one of each kind that succeeds. If the row is concurrently deleted, then `MERGE` will evaluate the command's `NOT MATCHED [BY TARGET]` actions, and execute the first one that succeeds. If `MERGE` attempts an `INSERT` and a unique index is present and a duplicate row is concurrently inserted, then a uniqueness violation error is raised; `MERGE` does not attempt to avoid such errors by restarting evaluation of `MATCHED` conditions.

Because of the above rules, it is possible for an updating command to see an inconsistent snapshot: it can see the effects of concurrent updating commands on the same rows it is trying to update, but it does not see effects of those commands on other rows in the database. This behavior makes Read Committed mode unsuitable for commands that involve complex search conditions; however, it is just right for simpler cases. For example, consider updating bank balances with transactions like:

```
BEGIN;
UPDATE accounts SET balance = balance + 100.00 WHERE acctnum = 12345;
UPDATE accounts SET balance = balance - 100.00 WHERE acctnum = 7534;
COMMIT;
```

If two such transactions concurrently try to change the balance of account 12345, we clearly want the second transaction to start with the updated version of the account's row. Because each command is affecting only a predetermined row, letting it see the updated version of the row does not create any troublesome inconsistency.

More complex usage can produce undesirable results in Read Committed mode. For example, consider a `DELETE` command operating on data that is being both added and removed from its restriction criteria by another command, e.g., assume `website` is a two-row table with `website.hits` equaling 9 and 10:

```
BEGIN;
UPDATE website SET hits = hits + 1;
-- run from another session: DELETE FROM website WHERE hits = 10;
COMMIT;
```

The `DELETE` will have no effect even though there is a `website.hits = 10` row before and after the `UPDATE`. This occurs because the pre-update row value 9 is skipped, and when the `UPDATE` completes and `DELETE` obtains a lock, the new row value is no longer 10 but 11, which no longer matches the criteria.

Because Read Committed mode starts each command with a new snapshot that includes all transactions committed up to that instant, subsequent commands in the same transaction will see the effects of the committed concurrent transaction in any case. The point at issue above is whether or not a *single* command sees an absolutely consistent view of the database.

The partial transaction isolation provided by Read Committed mode is adequate for many applications, and this mode is fast and simple to use; however, it is not sufficient for all cases. Applications that do complex queries and updates might require a more rigorously consistent view of the database than Read Committed mode provides.

13.2.2. Repeatable Read Isolation Level

The *Repeatable Read* isolation level only sees data committed before the transaction began; it never sees either uncommitted data or changes committed by concurrent transactions during the transaction's execution. (However, each query does see the effects of previous updates executed within its own transaction, even though they are not yet committed.) This is a stronger guarantee than is required by the SQL standard for this isolation level, and prevents all of the phenomena described in [Table 13.1](#) except for serialization anomalies. As mentioned above, this is specifically allowed by the standard, which only describes the *minimum* protections each isolation level must provide.

This level is different from Read Committed in that a query in a repeatable read transaction sees a snapshot as of the start of the first non-transaction-control statement in the *transaction*, not as of the start of the current statement within the transaction. Thus, successive `SELECT` commands within a *single* transaction see the same data, i.e., they do not see changes made by other transactions that committed after their own transaction started.

Applications using this level must be prepared to retry transactions due to serialization failures.

`UPDATE`, `DELETE`, `MERGE`, `SELECT FOR UPDATE`, and `SELECT FOR SHARE` commands behave the same as `SELECT` in terms of searching for target rows: they will only find target rows that were committed as of the transaction start time. However, such a target row might have already been updated (or deleted or locked) by another concurrent transaction by the time it is found. In this case, the repeatable read transaction will wait for the first updating transaction to commit or roll back (if it is still in progress). If the first updater rolls back, then its effects are negated and the repeatable read transaction can proceed with updating the originally found row. But if the first updater commits (and actually updated or deleted the row, not just locked it) then the repeatable read transaction will be rolled back with the message

```
ERROR:  could not serialize access due to concurrent update
```

because a repeatable read transaction cannot modify or lock rows changed by other transactions after the repeatable read transaction began.

When an application receives this error message, it should abort the current transaction and retry the whole transaction from the beginning. The second time through, the transaction will see the previously-committed change as part of its initial view of the database, so there is no logical conflict in using the new version of the row as the starting point for the new transaction's update.

Note that only updating transactions might need to be retried; read-only transactions will never have serialization conflicts.

The Repeatable Read mode provides a rigorous guarantee that each transaction sees a completely stable view of the database. However, this view will not necessarily always be consistent with some serial (one at a time) execution of concurrent transactions of the same level. For example, even a read-only transaction at this level may see a control record updated to show that a batch has been completed but *not* see one of the detail records which is logically part of the batch because it read an earlier revision of the control record. Attempts to enforce business rules by transactions running at this isolation level are not likely to work correctly without careful use of explicit locks to block conflicting transactions.

The Repeatable Read isolation level is implemented using a technique known in academic database literature and in some other database products as *Snapshot Isolation*. Differences in behavior and performance may be observed when compared with systems that use a traditional locking technique that reduces concurrency. Some other systems may even offer Repeatable Read and Snapshot Isolation as distinct isolation levels with different behavior. The permitted phenomena that distinguish the two techniques were not formalized by database researchers until after the SQL standard was developed, and are outside the scope of this manual. For a full treatment, please see [berenson95](#).

Note

Prior to PostgreSQL version 9.1, a request for the Serializable transaction isolation level provided exactly the same behavior described here. To retain the legacy Serializable behavior, Repeatable Read should now be requested.

13.2.3. Serializable Isolation Level

The *Serializable* isolation level provides the strictest transaction isolation. This level emulates serial transaction execution for all committed transactions; as if transactions had been executed one after another, serially, rather than concurrently. However, like the Repeatable Read level, applications using this level must be prepared to retry transactions due to serialization failures. In fact, this isolation level works exactly the same as Repeatable Read except that it also monitors for conditions which could make execution of a concurrent set of serializable transactions behave in a manner inconsistent with all possible serial (one at a time) executions of those transactions. This monitoring does not introduce any blocking beyond that present in repeatable read, but there is some overhead to the monitoring, and detection of the conditions which could cause a *serialization anomaly* will trigger a *serialization failure*.

As an example, consider a table `mytab`, initially containing:

class	value
1	10
1	20
2	100
2	200

Suppose that serializable transaction A computes:

```
SELECT SUM(value) FROM mytab WHERE class = 1;
```

and then inserts the result (30) as the `value` in a new row with `class = 2`. Concurrently, serializable transaction B computes:

```
SELECT SUM(value) FROM mytab WHERE class = 2;
```

and obtains the result 300, which it inserts in a new row with `class = 1`. Then both transactions try to commit. If either transaction were running at the Repeatable Read isolation level, both would be allowed to commit; but since there is no serial order of execution consistent with the result, using Serializable transactions will allow one transaction to commit and will roll the other back with this message:

```
ERROR: could not serialize access due to read/write dependencies among transactions
```

This is because if A had executed before B, B would have computed the sum 330, not 300, and similarly the other order would have resulted in a different sum computed by A.

When relying on Serializable transactions to prevent anomalies, it is important that any data read from a permanent user table not be considered valid until the transaction which read it has successfully committed. This is true even for read-only transactions, except that data read within a *deferrable* read-only transaction is known to be valid as soon as it is read, because such a transaction waits until it can acquire a snapshot guaranteed to be free from such problems before starting to read any data. In all other cases applications must not depend on results read during a transaction that later aborted; instead, they should retry the transaction until it succeeds.

To guarantee true serializability Postgres Pro uses *predicate locking*, which means that it keeps locks which allow it to determine when a write would have had an impact on the result of a previous read from a concurrent transaction, had it run first. In Postgres Pro these locks do not cause any blocking and therefore can *not* play any part in causing a deadlock. They are used to identify and flag dependencies among concurrent Serializable transactions which in certain combinations can lead to serialization anomalies. In contrast, a Read Committed or Repeatable Read transaction which wants to ensure data consistency may need to take out a lock on an entire table, which could block other users attempting to

use that table, or it may use `SELECT FOR UPDATE` or `SELECT FOR SHARE` which not only can block other transactions but cause disk access.

Predicate locks in Postgres Pro, like in most other database systems, are based on data actually accessed by a transaction. These will show up in the `pg_locks` system view with a mode of `SIReadLock`. The particular locks acquired during execution of a query will depend on the plan used by the query, and multiple finer-grained locks (e.g., tuple locks) may be combined into fewer coarser-grained locks (e.g., page locks) during the course of the transaction to prevent exhaustion of the memory used to track the locks. A `READ ONLY` transaction may be able to release its `SIRead` locks before completion, if it detects that no conflicts can still occur which could lead to a serialization anomaly. In fact, `READ ONLY` transactions will often be able to establish that fact at startup and avoid taking any predicate locks. If you explicitly request a `SERIALIZABLE READ ONLY DEFERRABLE` transaction, it will block until it can establish this fact. (This is the *only* case where Serializable transactions block but Repeatable Read transactions don't.) On the other hand, `SIRead` locks often need to be kept past transaction commit, until overlapping read write transactions complete.

Consistent use of Serializable transactions can simplify development. The guarantee that any set of successfully committed concurrent Serializable transactions will have the same effect as if they were run one at a time means that if you can demonstrate that a single transaction, as written, will do the right thing when run by itself, you can have confidence that it will do the right thing in any mix of Serializable transactions, even without any information about what those other transactions might do, or it will not successfully commit. It is important that an environment which uses this technique have a generalized way of handling serialization failures (which always return with an `SQLSTATE` value of `'40001'`), because it will be very hard to predict exactly which transactions might contribute to the read/write dependencies and need to be rolled back to prevent serialization anomalies. The monitoring of read/write dependencies has a cost, as does the restart of transactions which are terminated with a serialization failure, but balanced against the cost and blocking involved in use of explicit locks and `SELECT FOR UPDATE` or `SELECT FOR SHARE`, Serializable transactions are the best performance choice for some environments.

While Postgres Pro's Serializable transaction isolation level only allows concurrent transactions to commit if it can prove there is a serial order of execution that would produce the same effect, it doesn't always prevent errors from being raised that would not occur in true serial execution. In particular, it is possible to see unique constraint violations caused by conflicts with overlapping Serializable transactions even after explicitly checking that the key isn't present before attempting to insert it. This can be avoided by making sure that *all* Serializable transactions that insert potentially conflicting keys explicitly check if they can do so first. For example, imagine an application that asks the user for a new key and then checks that it doesn't exist already by trying to select it first, or generates a new key by selecting the maximum existing key and adding one. If some Serializable transactions insert new keys directly without following this protocol, unique constraints violations might be reported even in cases where they could not occur in a serial execution of the concurrent transactions.

For optimal performance when relying on Serializable transactions for concurrency control, these issues should be considered:

- Declare transactions as `READ ONLY` when possible.
- Control the number of active connections, using a connection pool if needed. This is always an important performance consideration, but it can be particularly important in a busy system using Serializable transactions.
- Don't put more into a single transaction than needed for integrity purposes.
- Don't leave connections dangling “idle in transaction” longer than necessary. The configuration parameter `idle_in_transaction_session_timeout` may be used to automatically disconnect lingering sessions.
- Eliminate explicit locks, `SELECT FOR UPDATE`, and `SELECT FOR SHARE` where no longer needed due to the protections automatically provided by Serializable transactions.
- When the system is forced to combine multiple page-level predicate locks into a single relation-level predicate lock because the predicate lock table is short of memory, an increase in the rate of se-

rialization failures may occur. You can avoid this by increasing `max_pred_locks_per_transaction`, `max_pred_locks_per_relation`, and/or `max_pred_locks_per_page`.

- A sequential scan will always necessitate a relation-level predicate lock. This can result in an increased rate of serialization failures. It may be helpful to encourage the use of index scans by reducing `random_page_cost` and/or increasing `cpu_tuple_cost`. Be sure to weigh any decrease in transaction rollbacks and restarts against any overall change in query execution time.

The Serializable isolation level is implemented using a technique known in academic database literature as Serializable Snapshot Isolation, which builds on Snapshot Isolation by adding checks for serialization anomalies. Some differences in behavior and performance may be observed when compared with other systems that use a traditional locking technique. Please see [ports12](#) for detailed information.

13.3. Explicit Locking

Postgres Pro provides various lock modes to control concurrent access to data in tables. These modes can be used for application-controlled locking in situations where MVCC does not give the desired behavior. Also, most Postgres Pro commands automatically acquire locks of appropriate modes to ensure that referenced tables are not dropped or modified in incompatible ways while the command executes. (For example, `TRUNCATE` cannot safely be executed concurrently with other operations on the same table, so it obtains an `ACCESS EXCLUSIVE` lock on the table to enforce that.)

To examine a list of the currently outstanding locks in a database server, use the `pg_locks` system view. For more information on monitoring the status of the lock manager subsystem, refer to [Chapter 27](#).

13.3.1. Table-Level Locks

The list below shows the available lock modes and the contexts in which they are used automatically by Postgres Pro. You can also acquire any of these locks explicitly with the command `LOCK`. Remember that all of these lock modes are table-level locks, even if the name contains the word “row”; the names of the lock modes are historical. To some extent the names reflect the typical usage of each lock mode — but the semantics are all the same. The only real difference between one lock mode and another is the set of lock modes with which each conflicts (see [Table 13.2](#)). Two transactions cannot hold locks of conflicting modes on the same table at the same time. (However, a transaction never conflicts with itself. For example, it might acquire `ACCESS EXCLUSIVE` lock and later acquire `ACCESS SHARE` lock on the same table.) Non-conflicting lock modes can be held concurrently by many transactions. Notice in particular that some lock modes are self-conflicting (for example, an `ACCESS EXCLUSIVE` lock cannot be held by more than one transaction at a time) while others are not self-conflicting (for example, an `ACCESS SHARE` lock can be held by multiple transactions).

Table-Level Lock Modes

`ACCESS SHARE` (`AccessShareLock`)

Conflicts with the `ACCESS EXCLUSIVE` lock mode only.

The `SELECT` command acquires a lock of this mode on referenced tables. In general, any query that only *reads* a table and does not modify it will acquire this lock mode.

`ROW SHARE` (`RowShareLock`)

Conflicts with the `EXCLUSIVE` and `ACCESS EXCLUSIVE` lock modes.

The `SELECT` command acquires a lock of this mode on all tables on which one of the `FOR UPDATE`, `FOR NO KEY UPDATE`, `FOR SHARE`, or `FOR KEY SHARE` options is specified (in addition to `ACCESS SHARE` locks on any other tables that are referenced without any explicit `FOR . . . locking` option).

`ROW EXCLUSIVE` (`RowExclusiveLock`)

Conflicts with the `SHARE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` lock modes.

The commands `UPDATE`, `DELETE`, `INSERT`, and `MERGE` acquire this lock mode on the target table (in addition to `ACCESS SHARE` locks on any other referenced tables). In general, this lock mode will be acquired by any command that *modifies data* in a table.

`SHARE UPDATE EXCLUSIVE` (`ShareUpdateExclusiveLock`)

Conflicts with the `SHARE UPDATE EXCLUSIVE`, `SHARE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` lock modes. This mode protects a table against concurrent schema changes and `VACUUM` runs.

Acquired by `VACUUM` (without `FULL`), `ANALYZE`, `CREATE INDEX CONCURRENTLY`, `CREATE STATISTICS`, `COMMENT ON`, `REINDEX CONCURRENTLY`, and certain `ALTER INDEX` and `ALTER TABLE` variants (for full details see the documentation of these commands).

`SHARE` (`ShareLock`)

Conflicts with the `ROW EXCLUSIVE`, `SHARE UPDATE EXCLUSIVE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` lock modes. This mode protects a table against concurrent data changes.

Acquired by `CREATE INDEX` (without `CONCURRENTLY`).

`SHARE ROW EXCLUSIVE` (`ShareRowExclusiveLock`)

Conflicts with the `ROW EXCLUSIVE`, `SHARE UPDATE EXCLUSIVE`, `SHARE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` lock modes. This mode protects a table against concurrent data changes, and is self-exclusive so that only one session can hold it at a time.

Acquired by `CREATE TRIGGER` and some forms of `ALTER TABLE`.

`EXCLUSIVE` (`ExclusiveLock`)

Conflicts with the `ROW SHARE`, `ROW EXCLUSIVE`, `SHARE UPDATE EXCLUSIVE`, `SHARE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` lock modes. This mode allows only concurrent `ACCESS SHARE` locks, i.e., only reads from the table can proceed in parallel with a transaction holding this lock mode.

Acquired by `REFRESH MATERIALIZED VIEW CONCURRENTLY`.

`ACCESS EXCLUSIVE` (`AccessExclusiveLock`)

Conflicts with locks of `ACCESS SHARE`, `ROW SHARE`, `ROW EXCLUSIVE`, `SHARE UPDATE EXCLUSIVE`, `SHARE`, `SHARE ROW EXCLUSIVE`, `EXCLUSIVE`, and `ACCESS EXCLUSIVE` mode. This mode guarantees that the holder is the only transaction accessing the table in any way.

Acquired by the `DROP TABLE`, `TRUNCATE`, `REINDEX`, `CLUSTER`, `VACUUM FULL`, and `REFRESH MATERIALIZED VIEW` (without `CONCURRENTLY`) commands. Many forms of `ALTER INDEX` and `ALTER TABLE` also acquire a lock at this level. This is also the default lock mode for `LOCK TABLE` statements that do not specify a mode explicitly.

Tip

Only an `ACCESS EXCLUSIVE` lock blocks a `SELECT` (without `FOR UPDATE/SHARE`) statement.

Besides, there are two additional lock modes designed to support 1C Enterprise. These modes do not conflict with any other lock modes described above. Their use is possible, but not encouraged, as [advisory locks](#) provide the same functionality.

`APPLICATION SHARE`

Conflicts with the `APPLICATION EXCLUSIVE` lock mode only.

APPLICATION EXCLUSIVE

Conflicts with the APPLICATION SHARE and APPLICATION EXCLUSIVE lock modes.

Once acquired, a lock is normally held until the end of the transaction. But if a lock is acquired after establishing a savepoint, the lock is released immediately if the savepoint is rolled back to. This is consistent with the principle that ROLLBACK cancels all effects of the commands since the savepoint. The same holds for locks acquired within a PL/pgSQL exception block: an error escape from the block releases locks acquired within it.

Table 13.2. Conflicting Lock Modes

Request- ed Lock Mode	Existing Lock Mode									
	ACCESS SHARE	ROW SHARE	ROW EX- CL.	SHARE UPDATE EXCL.	SHARE	SHARE ROW EX- CL.	EXCL.	ACCESS EXCL.	APP. SHARE	APP. EXCL.
ACCESS SHARE								X		
ROW SHARE							X	X		
ROW EX- CL.					X	X	X	X		
SHARE UPDATE EXCL.				X	X	X	X	X		
SHARE			X	X		X	X	X		
SHARE ROW EX- CL.			X	X	X	X	X	X		
EXCL.		X	X	X	X	X	X	X		
ACCESS EXCL.	X	X	X	X	X	X	X	X		
APP. SHARE										X
APP. EX- CL.									X	X

13.3.2. Row-Level Locks

In addition to table-level locks, there are row-level locks, which are listed as below with the contexts in which they are used automatically by Postgres Pro. See [Table 13.3](#) for a complete table of row-level lock conflicts. Note that a transaction can hold conflicting locks on the same row, even in different subtransactions; but other than that, two transactions can never hold conflicting locks on the same row. Row-level locks do not affect data querying; they block only *writers and lockers* to the same row. Row-level locks are released at transaction end or during savepoint rollback, just like table-level locks.

Row-Level Lock Modes

FOR UPDATE

FOR UPDATE causes the rows retrieved by the SELECT statement to be locked as though for update. This prevents them from being locked, modified or deleted by other transactions until the current transaction ends. That is, other transactions that attempt UPDATE, DELETE, SELECT FOR UPDATE, SELECT FOR NO KEY UPDATE, SELECT FOR SHARE or SELECT FOR KEY SHARE of these rows will be blocked until the current transaction ends; conversely, SELECT FOR UPDATE will wait for a concurrent

transaction that has run any of those commands on the same row, and will then lock and return the updated row (or no row, if the row was deleted). Within a `REPEATABLE READ` or `SERIALIZABLE` transaction, however, an error will be thrown if a row to be locked has changed since the transaction started. For further discussion see [Section 13.4](#).

The `FOR UPDATE` lock mode is also acquired by any `DELETE` on a row, and also by an `UPDATE` that modifies the values of certain columns. Currently, the set of columns considered for the `UPDATE` case are those that have a unique index on them that can be used in a foreign key (so partial indexes and expressional indexes are not considered), but this may change in the future.

`FOR NO KEY UPDATE`

Behaves similarly to `FOR UPDATE`, except that the lock acquired is weaker: this lock will not block `SELECT FOR KEY SHARE` commands that attempt to acquire a lock on the same rows. This lock mode is also acquired by any `UPDATE` that does not acquire a `FOR UPDATE` lock.

`FOR SHARE`

Behaves similarly to `FOR NO KEY UPDATE`, except that it acquires a shared lock rather than exclusive lock on each retrieved row. A shared lock blocks other transactions from performing `UPDATE`, `DELETE`, `SELECT FOR UPDATE` or `SELECT FOR NO KEY UPDATE` on these rows, but it does not prevent them from performing `SELECT FOR SHARE` or `SELECT FOR KEY SHARE`.

`FOR KEY SHARE`

Behaves similarly to `FOR SHARE`, except that the lock is weaker: `SELECT FOR UPDATE` is blocked, but not `SELECT FOR NO KEY UPDATE`. A key-shared lock blocks other transactions from performing `DELETE` or any `UPDATE` that changes the key values, but not other `UPDATE`, and neither does it prevent `SELECT FOR NO KEY UPDATE`, `SELECT FOR SHARE`, or `SELECT FOR KEY SHARE`.

Postgres Pro doesn't remember any information about modified rows in memory, so there is no limit on the number of rows locked at one time. However, locking a row might cause a disk write, e.g., `SELECT FOR UPDATE` modifies selected rows to mark them locked, and so will result in disk writes.

Table 13.3. Conflicting Row-Level Locks

Requested Lock Mode	Current Lock Mode			
	<code>FOR KEY SHARE</code>	<code>FOR SHARE</code>	<code>FOR NO KEY UPDATE</code>	<code>FOR UPDATE</code>
<code>FOR KEY SHARE</code>				X
<code>FOR SHARE</code>			X	X
<code>FOR NO KEY UPDATE</code>		X	X	X
<code>FOR UPDATE</code>	X	X	X	X

13.3.2.1. Specifics of `SELECT FOR UPDATE` in Distributed System

The behavior of the `SELECT FOR UPDATE` command in a distributed system may differ from the usual behavior in regular databases.

When used for a sharded table, `SELECT FOR UPDATE` may lock more rows than expected when rows are additionally filtered locally, `LIMIT` expressions are used, or multiple tables are joined. In these cases, entire tables may be locked despite the specified filters, which does not happen in non-distributed systems.

Example:

```
SELECT *
FROM bookings
WHERE book_date = bookings.now()::date
FOR UPDATE
```


This query locks 3 rows of the `booking` table on the local node and all rows of foreign table partitions on remote nodes. In the query execution plan, the remote query looks as follows:

```
...
Remote SQL: SELECT book_ref, book_date, total_amount, ctid FROM bookings.bookings_1 FOR
UPDATE
...
```

When running for a global table, `SELECT FOR UPDATE` locks only the corresponding local table. So, two users working on different nodes may lock the same rows simultaneously, which should not happen. This can lead to a deadlock, even if users make efforts to avoid it, for instance, by setting their own locks in advance or skipping other locks.

The first user runs the following query:

```
SELECT * FROM flights WHERE flight_id = 1976 FOR UPDATE
```

The query is executed successfully, and the row for the specified flight is locked.

The second user runs the same query and locks the same row on another node. Then, the second user tries to update the locked row:

```
UPDATE flights SET status = 'On Time' WHERE flight_id = 1976;
```

Now the second user encounters a lock. This query tries to update the corresponding rows on all remote nodes but one of them is already locked by the first user. Any further operations with these rows may result in a deadlock.

13.3.3. Page-Level Locks

In addition to table and row locks, page-level share/exclusive locks are used to control read/write access to table pages in the shared buffer pool. These locks are released immediately after a row is fetched or updated. Application developers normally need not be concerned with page-level locks, but they are mentioned here for completeness.

13.3.4. Deadlocks

The use of explicit locking can increase the likelihood of *deadlocks*, wherein two (or more) transactions each hold locks that the other wants. For example, if transaction 1 acquires an exclusive lock on table A and then tries to acquire an exclusive lock on table B, while transaction 2 has already exclusive-locked table B and now wants an exclusive lock on table A, then neither one can proceed. Postgres Pro automatically detects deadlock situations and resolves them by aborting one of the transactions involved, allowing the other(s) to complete. (Exactly which transaction will be aborted is difficult to predict and should not be relied upon.)

Note that deadlocks can also occur as the result of row-level locks (and thus, they can occur even if explicit locking is not used). Consider the case in which two concurrent transactions modify a table. The first transaction executes:

```
UPDATE accounts SET balance = balance + 100.00 WHERE acctnum = 11111;
```

This acquires a row-level lock on the row with the specified account number. Then, the second transaction executes:

```
UPDATE accounts SET balance = balance + 100.00 WHERE acctnum = 22222;
UPDATE accounts SET balance = balance - 100.00 WHERE acctnum = 11111;
```

The first `UPDATE` statement successfully acquires a row-level lock on the specified row, so it succeeds in updating that row. However, the second `UPDATE` statement finds that the row it is attempting to update has already been locked, so it waits for the transaction that acquired the lock to complete. Transaction two is now waiting on transaction one to complete before it continues execution. Now, transaction one executes:


```
UPDATE accounts SET balance = balance - 100.00 WHERE acctnum = 22222;
```

Transaction one attempts to acquire a row-level lock on the specified row, but it cannot: transaction two already holds such a lock. So it waits for transaction two to complete. Thus, transaction one is blocked on transaction two, and transaction two is blocked on transaction one: a deadlock condition. Postgres Pro will detect this situation and abort one of the transactions.

The best defense against deadlocks is generally to avoid them by being certain that all applications using a database acquire locks on multiple objects in a consistent order. In the example above, if both transactions had updated the rows in the same order, no deadlock would have occurred. One should also ensure that the first lock acquired on an object in a transaction is the most restrictive mode that will be needed for that object. If it is not feasible to verify this in advance, then deadlocks can be handled on-the-fly by retrying transactions that abort due to deadlocks.

So long as no deadlock situation is detected, a transaction seeking either a table-level or row-level lock will wait indefinitely for conflicting locks to be released. This means it is a bad idea for applications to hold transactions open for long periods of time (e.g., while waiting for user input).

13.3.5. Advisory Locks

Postgres Pro provides a means for creating locks that have application-defined meanings. These are called *advisory locks*, because the system does not enforce their use — it is up to the application to use them correctly. Advisory locks can be useful for locking strategies that are an awkward fit for the MVCC model. For example, a common use of advisory locks is to emulate pessimistic locking strategies typical of so-called “flat file” data management systems. While a flag stored in a table could be used for the same purpose, advisory locks are faster, avoid table bloat, and are automatically cleaned up by the server at the end of the session.

There are two ways to acquire an advisory lock in Postgres Pro: at session level or at transaction level. Once acquired at session level, an advisory lock is held until explicitly released or the session ends. Unlike standard lock requests, session-level advisory lock requests do not honor transaction semantics: a lock acquired during a transaction that is later rolled back will still be held following the rollback, and likewise an unlock is effective even if the calling transaction fails later. A lock can be acquired multiple times by its owning process; for each completed lock request there must be a corresponding unlock request before the lock is actually released. Transaction-level lock requests, on the other hand, behave more like regular lock requests: they are automatically released at the end of the transaction, and there is no explicit unlock operation. This behavior is often more convenient than the session-level behavior for short-term usage of an advisory lock. Session-level and transaction-level lock requests for the same advisory lock identifier will block each other in the expected way. If a session already holds a given advisory lock, additional requests by it will always succeed, even if other sessions are awaiting the lock; this statement is true regardless of whether the existing lock hold and new request are at session level or transaction level.

Like all locks in Postgres Pro, a complete list of advisory locks currently held by any session can be found in the [pg_locks](#) system view.

Both advisory locks and regular locks are stored in a shared memory pool whose size is defined by the configuration variables [max_locks_per_transaction](#) and [max_connections](#). Care must be taken not to exhaust this memory or the server will be unable to grant any locks at all. This imposes an upper limit on the number of advisory locks grantable by the server, typically in the tens to hundreds of thousands depending on how the server is configured.

In certain cases using advisory locking methods, especially in queries involving explicit ordering and `LIMIT` clauses, care must be taken to control the locks acquired because of the order in which SQL expressions are evaluated. For example:

```
SELECT pg_advisory_lock(id) FROM foo WHERE id = 12345; -- ok
SELECT pg_advisory_lock(id) FROM foo WHERE id > 12345 LIMIT 100; -- danger!
SELECT pg_advisory_lock(q.id) FROM
(
```

```
SELECT id FROM foo WHERE id > 12345 LIMIT 100
) q; -- ok
```

In the above queries, the second form is dangerous because the `LIMIT` is not guaranteed to be applied before the locking function is executed. This might cause some locks to be acquired that the application was not expecting, and hence would fail to release (until it ends the session). From the point of view of the application, such locks would be dangling, although still viewable in `pg_locks`.

The functions provided to manipulate advisory locks are described in [Section 9.28.10](#).

13.4. Data Consistency Checks at the Application Level

It is very difficult to enforce business rules regarding data integrity using Read Committed transactions because the view of the data is shifting with each statement, and even a single statement may not restrict itself to the statement's snapshot if a write conflict occurs.

While a Repeatable Read transaction has a stable view of the data throughout its execution, there is a subtle issue with using MVCC snapshots for data consistency checks, involving something known as *read/write conflicts*. If one transaction writes data and a concurrent transaction attempts to read the same data (whether before or after the write), it cannot see the work of the other transaction. The reader then appears to have executed first regardless of which started first or which committed first. If that is as far as it goes, there is no problem, but if the reader also writes data which is read by a concurrent transaction there is now a transaction which appears to have run before either of the previously mentioned transactions. If the transaction which appears to have executed last actually commits first, it is very easy for a cycle to appear in a graph of the order of execution of the transactions. When such a cycle appears, integrity checks will not work correctly without some help.

As mentioned in [Section 13.2.3](#), Serializable transactions are just Repeatable Read transactions which add nonblocking monitoring for dangerous patterns of read/write conflicts. When a pattern is detected which could cause a cycle in the apparent order of execution, one of the transactions involved is rolled back to break the cycle.

13.4.1. Enforcing Consistency with Serializable Transactions

If the Serializable transaction isolation level is used for all writes and for all reads which need a consistent view of the data, no other effort is required to ensure consistency. Software from other environments which is written to use serializable transactions to ensure consistency should “just work” in this regard in Postgres Pro.

When using this technique, it will avoid creating an unnecessary burden for application programmers if the application software goes through a framework which automatically retries transactions which are rolled back with a serialization failure. It may be a good idea to set `default_transaction_isolation` to `serializable`. It would also be wise to take some action to ensure that no other transaction isolation level is used, either inadvertently or to subvert integrity checks, through checks of the transaction isolation level in triggers.

See [Section 13.2.3](#) for performance suggestions.

Warning: Serializable Transactions and Data Replication

This level of integrity protection using Serializable transactions does not yet extend to hot standby mode ([Section 26.4](#)) or logical replicas. Because of that, those using hot standby or logical replication may want to use Repeatable Read and explicit locking on the primary.

13.4.2. Enforcing Consistency with Explicit Blocking Locks

When non-serializable writes are possible, to ensure the current validity of a row and protect it against concurrent updates one must use `SELECT FOR UPDATE`, `SELECT FOR SHARE`, or an appropriate `LOCK TABLE`

statement. (`SELECT FOR UPDATE` and `SELECT FOR SHARE` lock just the returned rows against concurrent updates, while `LOCK TABLE` locks the whole table.) This should be taken into account when porting applications to Postgres Pro from other environments.

Also of note to those converting from other environments is the fact that `SELECT FOR UPDATE` does not ensure that a concurrent transaction will not update or delete a selected row. To do that in Postgres Pro you must actually update the row, even if no values need to be changed. `SELECT FOR UPDATE` *temporarily blocks* other transactions from acquiring the same lock or executing an `UPDATE` or `DELETE` which would affect the locked row, but once the transaction holding this lock commits or rolls back, a blocked transaction will proceed with the conflicting operation unless an actual `UPDATE` of the row was performed while the lock was held.

Global validity checks require extra thought under non-serializable MVCC. For example, a banking application might wish to check that the sum of all credits in one table equals the sum of debits in another table, when both tables are being actively updated. Comparing the results of two successive `SELECT sum(...)` commands will not work reliably in Read Committed mode, since the second query will likely include the results of transactions not counted by the first. Doing the two sums in a single repeatable read transaction will give an accurate picture of only the effects of transactions that committed before the repeatable read transaction started — but one might legitimately wonder whether the answer is still relevant by the time it is delivered. If the repeatable read transaction itself applied some changes before trying to make the consistency check, the usefulness of the check becomes even more debatable, since now it includes some but not all post-transaction-start changes. In such cases a careful person might wish to lock all tables needed for the check, in order to get an indisputable picture of current reality. A `SHARE` mode (or higher) lock guarantees that there are no uncommitted changes in the locked table, other than those of the current transaction.

Note also that if one is relying on explicit locking to prevent concurrent changes, one should either use Read Committed mode, or in Repeatable Read mode be careful to obtain locks before performing queries. A lock obtained by a repeatable read transaction guarantees that no other transactions modifying the table are still running, but if the snapshot seen by the transaction predates obtaining the lock, it might predate some now-committed changes in the table. A repeatable read transaction's snapshot is actually frozen at the start of its first query or data-modification command (`SELECT`, `INSERT`, `UPDATE`, `DELETE`, or `MERGE`), so it is possible to obtain locks explicitly before the snapshot is frozen.

13.5. Serialization Failure Handling

Both Repeatable Read and Serializable isolation levels can produce errors that are designed to prevent serialization anomalies. As previously stated, applications using these levels must be prepared to retry transactions that fail due to serialization errors. Such an error's message text will vary according to the precise circumstances, but it will always have the SQLSTATE code 40001 (`serialization_failure`).

It may also be advisable to retry deadlock failures. These have the SQLSTATE code 40P01 (`deadlock_detected`).

In some cases it is also appropriate to retry unique-key failures, which have SQLSTATE code 23505 (`unique_violation`), and exclusion constraint failures, which have SQLSTATE code 23P01 (`exclusion_violation`). For example, if the application selects a new value for a primary key column after inspecting the currently stored keys, it could get a unique-key failure because another application instance selected the same new key concurrently. This is effectively a serialization failure, but the server will not detect it as such because it cannot “see” the connection between the inserted value and the previous reads. There are also some corner cases in which the server will issue a unique-key or exclusion constraint error even though in principle it has enough information to determine that a serialization problem is the underlying cause. While it's recommendable to just retry `serialization_failure` errors unconditionally, more care is needed when retrying these other error codes, since they might represent persistent error conditions rather than transient failures.

It is important to retry the complete transaction, including all logic that decides which SQL to issue and/or which values to use. Therefore, Postgres Pro does not offer an automatic retry facility, since it cannot do so with any guarantee of correctness.

Transaction retry does not guarantee that the retried transaction will complete; multiple retries may be needed. In cases with very high contention, it is possible that completion of a transaction may take many attempts. In cases involving a conflicting prepared transaction, it may not be possible to make progress until the prepared transaction commits or rolls back.

13.6. Caveats

Some DDL commands, currently only `TRUNCATE` and the table-rewriting forms of `ALTER TABLE`, are not MVCC-safe. This means that after the truncation or rewrite commits, the table will appear empty to concurrent transactions, if they are using a snapshot taken before the DDL command committed. This will only be an issue for a transaction that did not access the table in question before the DDL command started — any transaction that has done so would hold at least an `ACCESS SHARE` table lock, which would block the DDL command until that transaction completes. So these commands will not cause any apparent inconsistency in the table contents for successive queries on the target table, but they could cause visible inconsistency between the contents of the target table and other tables in the database.

Support for the Serializable transaction isolation level has not yet been added to hot standby replication targets (described in [Section 26.4](#)). The strictest isolation level currently supported in hot standby mode is Repeatable Read. While performing all permanent database writes within Serializable transactions on the primary will ensure that all standbys will eventually reach a consistent state, a Repeatable Read transaction run on the standby can sometimes see a transient state that is inconsistent with any serial execution of the transactions on the primary.

Internal access to the system catalogs is not done using the isolation level of the current transaction. This means that newly created database objects such as tables are visible to concurrent Repeatable Read and Serializable transactions, even though the rows they contain are not. In contrast, queries that explicitly examine the system catalogs don't see rows representing concurrently created database objects, in the higher isolation levels.

13.7. Locking and Indexes

Though Postgres Pro provides nonblocking read/write access to table data, nonblocking read/write access is not currently offered for every index access method implemented in Postgres Pro. The various index types are handled as follows:

B-tree, GiST and SP-GiST indexes

Short-term share/exclusive page-level locks are used for read/write access. Locks are released immediately after each index row is fetched or inserted. These index types provide the highest concurrency without deadlock conditions.

Hash indexes

Share/exclusive hash-bucket-level locks are used for read/write access. Locks are released after the whole bucket is processed. Bucket-level locks provide better concurrency than index-level ones, but deadlock is possible since the locks are held longer than one index operation.

GIN indexes

Short-term share/exclusive page-level locks are used for read/write access. Locks are released immediately after each index row is fetched or inserted. But note that insertion of a GIN-indexed value usually produces several index key insertions per row, so GIN might do substantial work for a single value's insertion.

Currently, B-tree indexes offer the best performance for concurrent applications; since they also have more features than hash indexes, they are the recommended index type for concurrent applications that need to index scalar data. When dealing with non-scalar data, B-trees are not useful, and GiST, SP-GiST or GIN indexes should be used instead.

Chapter 14. Performance Tips

Query performance can be affected by many things. Some of these can be controlled by the user, while others are fundamental to the underlying design of the system. This chapter provides some hints about understanding and tuning Postgres Pro performance.

Performance tuning should be done during application development and include an accurate choice of hardware (for example, estimating the number of CPUs and memory per Postgres Pro Shardman cluster node or tuning your storage), OS tuning (for example, tuning the `swappiness` parameter or network-related behavior) and DBMS tuning (choosing efficient configuration). But first of all, an application should be tested and tuned for distributed DBMS. This includes designing a distributed schema (or converting an existing schema to a distributed one), tuning queries, using connection poolers, caching and even checking performance issues related to possible serialization errors or Postgres Pro Shardman node outage. The design of the schema should include accurate selection of a sharding key and a decision which tables should become global. Usually you select a sharding key so that:

1. Most of the queries filter out most of sharded table partitions.
2. Sharded tables are colocated and all joins of sharded tables are equi-joins on the sharding key.

These rules allow Postgres Pro Shardman to efficiently exclude unused shards from queries and to push down joins to shards where the required data resides.

Each Postgres Pro Shardman node operates as a usual DBMS server, so all standard recommendations for tuning PostgreSQL for production load remain in place. You should select `shared_buffers`, `work_mem`, `effective_cache_size` depending on resources available to DBMS. Keep in mind that if the cluster topology is set to `cross`, Repfactor instances run on a single node `w`. When all cluster nodes are on-line, replicas should not utilize a lot of CPUs. However, in case of node failure, masters for Repfactor replication groups can become running on one server, which can create significant load on it. While tuning the `max_connections` parameter, note that each transaction can initiate `n-1` connections, where `n` is the number of replication groups in the cluster. When Silk is enabled, it is still true for transactions containing DML operations. When Silk is disabled, it is also true for read-only transactions.

Other parameters, which you perhaps would like to tune, are foreign server options. They can be set in `FDWOptions` section of Postgres Pro Shardman configuration file. Parameters that significantly affect Postgres Pro Shardman performance are `fetch_size`, `batch_size` and `async_capable`. When Silk transport is not enabled, `fetch_size` determines the number of records that are fetched from a remote server at once. When Silk transport is enabled, `fetch_size` currently does not have significant impact on the query execution. `batch_size` specifies how many rows can be combined in a single remote `INSERT` operation for a sharded table. `async_capable` allows asynchronous execution and should always be turned on (which is the default).

The `shardman.gt_batch_size` configuration parameter allows you to optimize the size of an intermediate buffer for `INSERT` and `DELETE` operations on global tables.

14.1. Using EXPLAIN

Postgres Pro devises a *query plan* for each query it receives. Choosing the right plan to match the query structure and the properties of the data is absolutely critical for good performance, so the system includes a complex *planner* that tries to choose good plans. You can use the `EXPLAIN` command to see what query plan the planner creates for any query. Plan-reading is an art that requires some experience to master, but this section attempts to cover the basics.

Examples in this section are drawn from the regression test database after doing a `VACUUM ANALYZE`, using v17 development sources. You should be able to get similar results if you try the examples yourself, but your estimated costs and row counts might vary slightly because `ANALYZE`'s statistics are random samples rather than exact, and because costs are inherently somewhat platform-dependent.

The examples use `EXPLAIN`'s default “text” output format, which is compact and convenient for humans to read. If you want to feed `EXPLAIN`'s output to a program for further analysis, you should use one of its machine-readable output formats (XML, JSON, or YAML) instead.

14.1.1. EXPLAIN Basics

The structure of a query plan is a tree of *plan nodes*. Nodes at the bottom level of the tree are scan nodes: they return raw rows from a table. There are different types of scan nodes for different table access methods: sequential scans, index scans, and bitmap index scans. There are also non-table row sources, such as `VALUES` clauses and set-returning functions in `FROM`, which have their own scan node types. If the query requires joining, aggregation, sorting, or other operations on the raw rows, then there will be additional nodes above the scan nodes to perform these operations. Again, there is usually more than one possible way to do these operations, so different node types can appear here too. The output of `EXPLAIN` has one line for each node in the plan tree, showing the basic node type plus the cost estimates that the planner made for the execution of that plan node. Additional lines might appear, indented from the node's summary line, to show additional properties of the node. The very first line (the summary line for the topmost node) has the estimated total execution cost for the plan; it is this number that the planner seeks to minimize.

Here is a trivial example, just to show what the output looks like. Note this applies for regular systems, for sharded ones refer to [Examining Plans](#).

```
EXPLAIN SELECT * FROM tenk1;
```

QUERY PLAN

```
-----
Seq Scan on tenk1  (cost=0.00..445.00 rows=10000 width=244)
```

Since this query has no `WHERE` clause, it must scan all the rows of the table, so the planner has chosen to use a simple sequential scan plan. The numbers that are quoted in parentheses are (left to right):

- Estimated start-up cost. This is the time expended before the output phase can begin, e.g., time to do the sorting in a sort node.
- Estimated total cost. This is stated on the assumption that the plan node is run to completion, i.e., all available rows are retrieved. In practice a node's parent node might stop short of reading all available rows (see the `LIMIT` example below).
- Estimated number of rows output by this plan node. Again, the node is assumed to be run to completion.
- Estimated average width of rows output by this plan node (in bytes).

The costs are measured in arbitrary units determined by the planner's cost parameters (see [Section 19.7.2](#)). Traditional practice is to measure the costs in units of disk page fetches; that is, `seq_page_cost` is conventionally set to 1.0 and the other cost parameters are set relative to that. The examples in this section are run with the default cost parameters.

It's important to understand that the cost of an upper-level node includes the cost of all its child nodes. It's also important to realize that the cost only reflects things that the planner cares about. In particular, the cost does not consider the time spent to convert output values to text form or to transmit them to the client, which could be important factors in the real elapsed time; but the planner ignores those costs because it cannot change them by altering the plan. (Every correct plan will output the same row set, we trust.)

The `rows` value is a little tricky because it is not the number of rows processed or scanned by the plan node, but rather the number emitted by the node. This is often less than the number scanned, as a result of filtering by any `WHERE`-clause conditions that are being applied at the node. Ideally the top-level rows estimate will approximate the number of rows actually returned, updated, or deleted by the query.

Returning to our example:

```
EXPLAIN SELECT * FROM tenk1;
```

QUERY PLAN

```
-----
Seq Scan on tenk1  (cost=0.00..445.00 rows=10000 width=244)
```


These numbers are derived very straightforwardly. If you do:

```
SELECT relpages, reltuples FROM pg_class WHERE relname = 'tenk1';
```

you will find that `tenk1` has 345 disk pages and 10000 rows. The estimated cost is computed as (disk pages read * [seq_page_cost](#)) + (rows scanned * [cpu_tuple_cost](#)). By default, `seq_page_cost` is 1.0 and `cpu_tuple_cost` is 0.01, so the estimated cost is $(345 * 1.0) + (10000 * 0.01) = 445$.

Now let's modify the query to add a `WHERE` condition:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 7000;
```

QUERY PLAN

```
-----
Seq Scan on tenk1  (cost=0.00..470.00 rows=7000 width=244)
  Filter: (unique1 < 7000)
```

Notice that the `EXPLAIN` output shows the `WHERE` clause being applied as a “filter” condition attached to the Seq Scan plan node. This means that the plan node checks the condition for each row it scans, and outputs only the ones that pass the condition. The estimate of output rows has been reduced because of the `WHERE` clause. However, the scan will still have to visit all 10000 rows, so the cost hasn't decreased; in fact it has gone up a bit (by $10000 * \text{cpu_operator_cost}$, to be exact) to reflect the extra CPU time spent checking the `WHERE` condition.

The actual number of rows this query would select is 7000, but the `rows` estimate is only approximate. If you try to duplicate this experiment, you may well get a slightly different estimate; moreover, it can change after each `ANALYZE` command, because the statistics produced by `ANALYZE` are taken from a randomized sample of the table.

Now, let's make the condition more restrictive:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 100;
```

QUERY PLAN

```
-----
Bitmap Heap Scan on tenk1  (cost=5.06..224.98 rows=100 width=244)
  Recheck Cond: (unique1 < 100)
    -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100 width=0)
        Index Cond: (unique1 < 100)
```

Here the planner has decided to use a two-step plan: the child plan node visits an index to find the locations of rows matching the index condition, and then the upper plan node actually fetches those rows from the table itself. Fetching rows separately is much more expensive than reading them sequentially, but because not all the pages of the table have to be visited, this is still cheaper than a sequential scan. (The reason for using two plan levels is that the upper plan node sorts the row locations identified by the index into physical order before reading them, to minimize the cost of separate fetches. The “bitmap” mentioned in the node names is the mechanism that does the sorting.)

Now let's add another condition to the `WHERE` clause:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 100 AND stringu1 = 'xxx';
```

QUERY PLAN

```
-----
Bitmap Heap Scan on tenk1  (cost=5.04..225.20 rows=1 width=244)
  Recheck Cond: (unique1 < 100)
  Filter: (stringu1 = 'xxx'::name)
    -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100 width=0)
        Index Cond: (unique1 < 100)
```

The added condition `stringu1 = 'xxx'` reduces the output row count estimate, but not the cost because we still have to visit the same set of rows. That's because the `stringu1` clause cannot be applied as an

index condition, since this index is only on the `unique1` column. Instead it is applied as a filter on the rows retrieved using the index. Thus the cost has actually gone up slightly to reflect this extra checking.

In some cases the planner will prefer a “simple” index scan plan:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 = 42;
```

QUERY PLAN

```
-----
Index Scan using tenk1_unique1 on tenk1  (cost=0.29..8.30 rows=1 width=244)
  Index Cond: (unique1 = 42)
```

In this type of plan the table rows are fetched in index order, which makes them even more expensive to read, but there are so few that the extra cost of sorting the row locations is not worth it. You'll most often see this plan type for queries that fetch just a single row. It's also often used for queries that have an `ORDER BY` condition that matches the index order, because then no extra sorting step is needed to satisfy the `ORDER BY`. In this example, adding `ORDER BY unique1` would use the same plan because the index already implicitly provides the requested ordering.

The planner may implement an `ORDER BY` clause in several ways. The above example shows that such an ordering clause may be implemented implicitly. The planner may also add an explicit `Sort` step:

```
EXPLAIN SELECT * FROM tenk1 ORDER BY unique1;
```

QUERY PLAN

```
-----
Sort  (cost=1109.39..1134.39 rows=10000 width=244)
  Sort Key: unique1
  -> Seq Scan on tenk1  (cost=0.00..445.00 rows=10000 width=244)
```

If a part of the plan guarantees an ordering on a prefix of the required sort keys, then the planner may instead decide to use an `Incremental Sort` step:

```
EXPLAIN SELECT * FROM tenk1 ORDER BY hundred, ten LIMIT 100;
```

QUERY PLAN

```
-----
Limit  (cost=19.35..39.49 rows=100 width=244)
  -> Incremental Sort  (cost=19.35..2033.39 rows=10000 width=244)
      Sort Key: hundred, ten
      Presorted Key: hundred
      -> Index Scan using tenk1_hundred on tenk1  (cost=0.29..1574.20 rows=10000
width=244)
```

Compared to regular sorts, sorting incrementally allows returning tuples before the entire result set has been sorted, which particularly enables optimizations with `LIMIT` queries. It may also reduce memory usage and the likelihood of spilling sorts to disk, but it comes at the cost of the increased overhead of splitting the result set into multiple sorting batches.

If there are separate indexes on several of the columns referenced in `WHERE`, the planner might choose to use an `AND` or `OR` combination of the indexes:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 100 AND unique2 > 9000;
```

QUERY PLAN

```
-----
Bitmap Heap Scan on tenk1  (cost=25.07..60.11 rows=10 width=244)
  Recheck Cond: ((unique1 < 100) AND (unique2 > 9000))
  -> BitmapAnd  (cost=25.07..25.07 rows=10 width=0)
      -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100 width=0)
          Index Cond: (unique1 < 100)
```



```
-> Bitmap Index Scan on tenk1_unique2 (cost=0.00..19.78 rows=999 width=0)
    Index Cond: (unique2 > 9000)
```

But this requires visiting both indexes, so it's not necessarily a win compared to using just one index and treating the other condition as a filter. If you vary the ranges involved you'll see the plan change accordingly.

Here is an example showing the effects of `LIMIT`:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 100 AND unique2 > 9000 LIMIT 2;
```

QUERY PLAN

```
-----
Limit (cost=0.29..14.28 rows=2 width=244)
-> Index Scan using tenk1_unique2 on tenk1 (cost=0.29..70.27 rows=10 width=244)
    Index Cond: (unique2 > 9000)
    Filter: (unique1 < 100)
```

This is the same query as above, but we added a `LIMIT` so that not all the rows need be retrieved, and the planner changed its mind about what to do. Notice that the total cost and row count of the Index Scan node are shown as if it were run to completion. However, the Limit node is expected to stop after retrieving only a fifth of those rows, so its total cost is only a fifth as much, and that's the actual estimated cost of the query. This plan is preferred over adding a Limit node to the previous plan because the Limit could not avoid paying the startup cost of the bitmap scan, so the total cost would be something over 25 units with that approach.

Let's try joining two tables, using the columns we have been discussing:

```
EXPLAIN SELECT *
FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 10 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```
-----
Nested Loop (cost=4.65..118.50 rows=10 width=488)
-> Bitmap Heap Scan on tenk1 t1 (cost=4.36..39.38 rows=10 width=244)
    Recheck Cond: (unique1 < 10)
-> Bitmap Index Scan on tenk1_unique1 (cost=0.00..4.36 rows=10 width=0)
    Index Cond: (unique1 < 10)
-> Index Scan using tenk2_unique2 on tenk2 t2 (cost=0.29..7.90 rows=1 width=244)
    Index Cond: (unique2 = t1.unique2)
```

In this plan, we have a nested-loop join node with two table scans as inputs, or children. The indentation of the node summary lines reflects the plan tree structure. The join's first, or "outer", child is a bitmap scan similar to those we saw before. Its cost and row count are the same as we'd get from `SELECT ... WHERE unique1 < 10` because we are applying the `WHERE` clause `unique1 < 10` at that node. The `t1.unique2 = t2.unique2` clause is not relevant yet, so it doesn't affect the row count of the outer scan. The nested-loop join node will run its second, or "inner" child once for each row obtained from the outer child. Column values from the current outer row can be plugged into the inner scan; here, the `t1.unique2` value from the outer row is available, so we get a plan and costs similar to what we saw above for a simple `SELECT ... WHERE t2.unique2 = constant` case. (The estimated cost is actually a bit lower than what was seen above, as a result of caching that's expected to occur during the repeated index scans on `t2`.) The costs of the loop node are then set on the basis of the cost of the outer scan, plus one repetition of the inner scan for each outer row ($10 * 7.90$, here), plus a little CPU time for join processing.

In this example the join's output row count is the same as the product of the two scans' row counts, but that's not true in all cases because there can be additional `WHERE` clauses that mention both tables and so can only be applied at the join point, not to either input scan. Here's an example:

```
EXPLAIN SELECT *
```

```
FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 10 AND t2.unique2 < 10 AND t1.hundred < t2.hundred;
```

QUERY PLAN

```
-----
Nested Loop  (cost=4.65..49.36 rows=33 width=488)
  Join Filter: (t1.hundred < t2.hundred)
  -> Bitmap Heap Scan on tenk1 t1  (cost=4.36..39.38 rows=10 width=244)
      Recheck Cond: (unique1 < 10)
      -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..4.36 rows=10 width=0)
          Index Cond: (unique1 < 10)
  -> Materialize  (cost=0.29..8.51 rows=10 width=244)
      -> Index Scan using tenk2_unique2 on tenk2 t2  (cost=0.29..8.46 rows=10
width=244)
          Index Cond: (unique2 < 10)
```

The condition `t1.hundred < t2.hundred` can't be tested in the `tenk2_unique2` index, so it's applied at the join node. This reduces the estimated output row count of the join node, but does not change either input scan.

Notice that here the planner has chosen to “materialize” the inner relation of the join, by putting a `Materialize` plan node atop it. This means that the `t2` index scan will be done just once, even though the nested-loop join node needs to read that data ten times, once for each row from the outer relation. The `Materialize` node saves the data in memory as it's read, and then returns the data from memory on each subsequent pass.

When dealing with outer joins, you might see join plan nodes with both “Join Filter” and plain “Filter” conditions attached. Join Filter conditions come from the outer join's `ON` clause, so a row that fails the Join Filter condition could still get emitted as a null-extended row. But a plain Filter condition is applied after the outer-join rules and so acts to remove rows unconditionally. In an inner join there is no semantic difference between these types of filters.

If we change the query's selectivity a bit, we might get a very different join plan:

```
EXPLAIN SELECT *
FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 100 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```
-----
Hash Join  (cost=226.23..709.73 rows=100 width=488)
  Hash Cond: (t2.unique2 = t1.unique2)
  -> Seq Scan on tenk2 t2  (cost=0.00..445.00 rows=10000 width=244)
  -> Hash  (cost=224.98..224.98 rows=100 width=244)
      -> Bitmap Heap Scan on tenk1 t1  (cost=5.06..224.98 rows=100 width=244)
          Recheck Cond: (unique1 < 100)
          -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100
width=0)
              Index Cond: (unique1 < 100)
```

Here, the planner has chosen to use a hash join, in which rows of one table are entered into an in-memory hash table, after which the other table is scanned and the hash table is probed for matches to each row. Again note how the indentation reflects the plan structure: the bitmap scan on `tenk1` is the input to the Hash node, which constructs the hash table. That's then returned to the Hash Join node, which reads rows from its outer child plan and searches the hash table for each one.

Another possible type of join is a merge join, illustrated here:

```
EXPLAIN SELECT *
```

```
FROM tenk1 t1, onek t2
WHERE t1.unique1 < 100 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```
-----
Merge Join  (cost=0.56..233.49 rows=10 width=488)
  Merge Cond: (t1.unique2 = t2.unique2)
    -> Index Scan using tenk1_unique2 on tenk1 t1  (cost=0.29..643.28 rows=100
width=244)
      Filter: (unique1 < 100)
    -> Index Scan using onek_unique2 on onek t2  (cost=0.28..166.28 rows=1000
width=244)
```

Merge join requires its input data to be sorted on the join keys. In this example each input is sorted by using an index scan to visit the rows in the correct order; but a sequential scan and sort could also be used. (Sequential-scan-and-sort frequently beats an index scan for sorting many rows, because of the nonsequential disk access required by the index scan.)

One way to look at variant plans is to force the planner to disregard whatever strategy it thought was the cheapest, using the enable/disable flags described in [Section 19.7.1](#). (This is a crude tool, but useful. See also [Section 14.3](#).) For example, if we're unconvinced that merge join is the best join type for the previous example, we could try

```
SET enable_mergejoin = off;
```

```
EXPLAIN SELECT *
FROM tenk1 t1, onek t2
WHERE t1.unique1 < 100 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```
-----
Hash Join  (cost=226.23..344.08 rows=10 width=488)
  Hash Cond: (t2.unique2 = t1.unique2)
    -> Seq Scan on onek t2  (cost=0.00..114.00 rows=1000 width=244)
    -> Hash  (cost=224.98..224.98 rows=100 width=244)
      -> Bitmap Heap Scan on tenk1 t1  (cost=5.06..224.98 rows=100 width=244)
        Recheck Cond: (unique1 < 100)
        -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100
width=0)
          Index Cond: (unique1 < 100)
```

which shows that the planner thinks that hash join would be nearly 50% more expensive than merge join for this case. Of course, the next question is whether it's right about that. We can investigate that using `EXPLAIN ANALYZE`, as discussed [below](#).

Some query plans involve *subplans*, which arise from sub-`SELECT`s in the original query. Such queries can sometimes be transformed into ordinary join plans, but when they cannot be, we get plans like:

```
EXPLAIN VERBOSE SELECT unique1
FROM tenk1 t
WHERE t.ten < ALL (SELECT o.ten FROM onek o WHERE o.four = t.four);
```

QUERY PLAN

```
-----
Seq Scan on public.tenk1 t  (cost=0.00..586095.00 rows=5000 width=4)
  Output: t.unique1
  Filter: (ALL (t.ten < (SubPlan 1).col1))
  SubPlan 1
```

```
-> Seq Scan on public.onek o  (cost=0.00..116.50 rows=250 width=4)
    Output: o.ten
    Filter: (o.four = t.four)
```

This rather artificial example serves to illustrate a couple of points: values from the outer plan level can be passed down into a subplan (here, `t.four` is passed down) and the results of the sub-select are available to the outer plan. Those result values are shown by `EXPLAIN` with notations like `(subplan_name).colN`, which refers to the *N*'th output column of the sub-`SELECT`.

In the example above, the `ALL` operator runs the subplan again for each row of the outer query (which accounts for the high estimated cost). Some queries can use a *hashed subplan* to avoid that:

```
EXPLAIN SELECT *
FROM tenk1 t
WHERE t.unique1 NOT IN (SELECT o.unique1 FROM onek o);
```

QUERY PLAN

```
-----
Seq Scan on tenk1 t  (cost=61.77..531.77 rows=5000 width=244)
    Filter: (NOT (ANY (unique1 = (hashed SubPlan 1).col1)))
    SubPlan 1
        -> Index Only Scan using onek_unique1 on onek o  (cost=0.28..59.27 rows=1000
            width=4)
        (4 rows)
```

Here, the subplan is run a single time and its output is loaded into an in-memory hash table, which is then probed by the outer `ANY` operator. This requires that the sub-`SELECT` not reference any variables of the outer query, and that the `ANY`'s comparison operator be amenable to hashing.

If, in addition to not referencing any variables of the outer query, the sub-`SELECT` cannot return more than one row, it may instead be implemented as an *initplan*:

```
EXPLAIN VERBOSE SELECT unique1
FROM tenk1 t1 WHERE t1.ten = (SELECT (random() * 10)::integer);
```

QUERY PLAN

```
-----
Seq Scan on public.tenk1 t1  (cost=0.02..470.02 rows=1000 width=4)
    Output: t1.unique1
    Filter: (t1.ten = (InitPlan 1).col1)
    InitPlan 1
        -> Result  (cost=0.00..0.02 rows=1 width=4)
            Output: ((random() * '10'::double precision))::integer
```

An *initplan* is run only once per execution of the outer plan, and its results are saved for re-use in later rows of the outer plan. So in this example `random()` is evaluated only once and all the values of `t1.ten` are compared to the same randomly-chosen integer. That's quite different from what would happen without the sub-`SELECT` construct.

14.1.2. EXPLAIN ANALYZE

It is possible to check the accuracy of the planner's estimates by using `EXPLAIN`'s `ANALYZE` option. With this option, `EXPLAIN` actually executes the query, and then displays the true row counts and true run time accumulated within each plan node, along with the same estimates that a plain `EXPLAIN` shows. For example, we might get a result like this:

```
EXPLAIN ANALYZE SELECT *
FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 10 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```

-----
Nested Loop  (cost=4.65..118.50 rows=10 width=488) (actual time=0.017..0.051 rows=10
loops=1)
-> Bitmap Heap Scan on tenk1 t1  (cost=4.36..39.38 rows=10 width=244) (actual
time=0.009..0.017 rows=10 loops=1)
    Recheck Cond: (unique1 < 10)
    Heap Blocks: exact=10
-> Bitmap Index Scan on tenk1_unique1  (cost=0.00..4.36 rows=10 width=0)
(actual time=0.004..0.004 rows=10 loops=1)
    Index Cond: (unique1 < 10)
-> Index Scan using tenk2_unique2 on tenk2 t2  (cost=0.29..7.90 rows=1 width=244)
(actual time=0.003..0.003 rows=1 loops=10)
    Index Cond: (unique2 = t1.unique2)
Planning Time: 0.485 ms
Execution Time: 0.073 ms

```

Note that the “actual time” values are in milliseconds of real time, whereas the `cost` estimates are expressed in arbitrary units; so they are unlikely to match up. The thing that's usually most important to look for is whether the estimated row counts are reasonably close to reality. In this example the estimates were all dead-on, but that's quite unusual in practice.

In some query plans, it is possible for a subplan node to be executed more than once. For example, the inner index scan will be executed once per outer row in the above nested-loop plan. In such cases, the `loops` value reports the total number of executions of the node, and the actual time and rows values shown are averages per-execution. This is done to make the numbers comparable with the way that the cost estimates are shown. Multiply by the `loops` value to get the total time actually spent in the node. In the above example, we spent a total of 0.030 milliseconds executing the index scans on `tenk2`.

In some cases `EXPLAIN ANALYZE` shows additional execution statistics beyond the plan node execution times and row counts. For example, Sort and Hash nodes provide extra information:

```

EXPLAIN ANALYZE SELECT *
FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 100 AND t1.unique2 = t2.unique2 ORDER BY t1.fivethous;

```

QUERY PLAN

```

-----
Sort  (cost=713.05..713.30 rows=100 width=488) (actual time=2.995..3.002 rows=100
loops=1)
    Sort Key: t1.fivethous
    Sort Method: quicksort  Memory: 74kB
-> Hash Join  (cost=226.23..709.73 rows=100 width=488) (actual time=0.515..2.920
rows=100 loops=1)
    Hash Cond: (t2.unique2 = t1.unique2)
-> Seq Scan on tenk2 t2  (cost=0.00..445.00 rows=10000 width=244) (actual
time=0.026..1.790 rows=10000 loops=1)
-> Hash  (cost=224.98..224.98 rows=100 width=244) (actual time=0.476..0.477
rows=100 loops=1)
    Buckets: 1024  Batches: 1  Memory Usage: 35kB
-> Bitmap Heap Scan on tenk1 t1  (cost=5.06..224.98 rows=100 width=244)
(actual time=0.030..0.450 rows=100 loops=1)
    Recheck Cond: (unique1 < 100)
    Heap Blocks: exact=90
-> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100
width=0) (actual time=0.013..0.013 rows=100 loops=1)
    Index Cond: (unique1 < 100)
Planning Time: 0.187 ms

```

Execution Time: 3.036 ms

The Sort node shows the sort method used (in particular, whether the sort was in-memory or on-disk) and the amount of memory or disk space needed. The Hash node shows the number of hash buckets and batches as well as the peak amount of memory used for the hash table. (If the number of batches exceeds one, there will also be disk space usage involved, but that is not shown.)

Another type of extra information is the number of rows removed by a filter condition:

```
EXPLAIN ANALYZE SELECT * FROM tenk1 WHERE ten < 7;
```

QUERY PLAN

```
-----  
Seq Scan on tenk1 (cost=0.00..470.00 rows=7000 width=244) (actual time=0.030..1.995  
rows=7000 loops=1)  
  Filter: (ten < 7)  
  Rows Removed by Filter: 3000  
Planning Time: 0.102 ms  
Execution Time: 2.145 ms
```

These counts can be particularly valuable for filter conditions applied at join nodes. The “Rows Removed” line only appears when at least one scanned row, or potential join pair in the case of a join node, is rejected by the filter condition.

A case similar to filter conditions occurs with “lossy” index scans. For example, consider this search for polygons containing a specific point:

```
EXPLAIN ANALYZE SELECT * FROM polygon_tbl WHERE f1 @> polygon '(0.5,2.0)';
```

QUERY PLAN

```
-----  
Seq Scan on polygon_tbl (cost=0.00..1.09 rows=1 width=85) (actual time=0.023..0.023  
rows=0 loops=1)  
  Filter: (f1 @> '((0.5,2))'::polygon)  
  Rows Removed by Filter: 7  
Planning Time: 0.039 ms  
Execution Time: 0.033 ms
```

The planner thinks (quite correctly) that this sample table is too small to bother with an index scan, so we have a plain sequential scan in which all the rows got rejected by the filter condition. But if we force an index scan to be used, we see:

```
SET enable_seqscan TO off;
```

```
EXPLAIN ANALYZE SELECT * FROM polygon_tbl WHERE f1 @> polygon '(0.5,2.0)';
```

QUERY PLAN

```
-----  
Index Scan using gpolygonind on polygon_tbl (cost=0.13..8.15 rows=1 width=85) (actual  
time=0.074..0.074 rows=0 loops=1)  
  Index Cond: (f1 @> '((0.5,2))'::polygon)  
  Rows Removed by Index Recheck: 1  
Planning Time: 0.039 ms  
Execution Time: 0.098 ms
```

Here we can see that the index returned one candidate row, which was then rejected by a recheck of the index condition. This happens because a GiST index is “lossy” for polygon containment tests: it actually returns the rows with polygons that overlap the target, and then we have to do the exact containment test on those rows.

EXPLAIN has a **BUFFERS** option that can be used with **ANALYZE** to get even more run time statistics:

```
EXPLAIN (ANALYZE, BUFFERS) SELECT * FROM tenk1 WHERE unique1 < 100 AND unique2 > 9000;
```

QUERY PLAN

```
-----
Bitmap Heap Scan on tenk1  (cost=25.07..60.11 rows=10 width=244) (actual
time=0.105..0.114 rows=10 loops=1)
  Recheck Cond: ((unique1 < 100) AND (unique2 > 9000))
  Heap Blocks: exact=10
  Buffers: shared hit=14 read=3
  -> BitmapAnd  (cost=25.07..25.07 rows=10 width=0) (actual time=0.100..0.101 rows=0
loops=1)
    Buffers: shared hit=4 read=3
    -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100 width=0)
(actual time=0.027..0.027 rows=100 loops=1)
      Index Cond: (unique1 < 100)
      Buffers: shared hit=2
    -> Bitmap Index Scan on tenk1_unique2  (cost=0.00..19.78 rows=999 width=0)
(actual time=0.070..0.070 rows=999 loops=1)
      Index Cond: (unique2 > 9000)
      Buffers: shared hit=2 read=3
Planning:
  Buffers: shared hit=3
Planning Time: 0.162 ms
Execution Time: 0.143 ms
```

The numbers provided by **BUFFERS** help to identify which parts of the query are the most I/O-intensive.

Keep in mind that because **EXPLAIN ANALYZE** actually runs the query, any side-effects will happen as usual, even though whatever results the query might output are discarded in favor of printing the **EXPLAIN** data. If you want to analyze a data-modifying query without changing your tables, you can roll the command back afterwards, for example:

```
BEGIN;
```

```
EXPLAIN ANALYZE UPDATE tenk1 SET hundred = hundred + 1 WHERE unique1 < 100;
```

QUERY PLAN

```
-----
Update on tenk1  (cost=5.06..225.23 rows=0 width=0) (actual time=1.634..1.635 rows=0
loops=1)
  -> Bitmap Heap Scan on tenk1  (cost=5.06..225.23 rows=100 width=10) (actual
time=0.065..0.141 rows=100 loops=1)
    Recheck Cond: (unique1 < 100)
    Heap Blocks: exact=90
    -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..5.04 rows=100 width=0)
(actual time=0.031..0.031 rows=100 loops=1)
      Index Cond: (unique1 < 100)
Planning Time: 0.151 ms
Execution Time: 1.856 ms
```

```
ROLLBACK;
```

As seen in this example, when the query is an **INSERT**, **UPDATE**, **DELETE**, or **MERGE** command, the actual work of applying the table changes is done by a top-level Insert, Update, Delete, or Merge plan node. The plan nodes underneath this node perform the work of locating the old rows and/or computing the new data. So above, we see the same sort of bitmap table scan we've seen already, and its output is

fed to an Update node that stores the updated rows. It's worth noting that although the data-modifying node can take a considerable amount of run time (here, it's consuming the lion's share of the time), the planner does not currently add anything to the cost estimates to account for that work. That's because the work to be done is the same for every correct query plan, so it doesn't affect planning decisions.

When an UPDATE, DELETE, or MERGE command affects a partitioned table or inheritance hierarchy, the output might look like this:

```
EXPLAIN UPDATE gtest_parent SET f1 = CURRENT_DATE WHERE f2 = 101;
```

QUERY PLAN

```
-----
Update on gtest_parent  (cost=0.00..3.06 rows=0 width=0)
  Update on gtest_child gtest_parent_1
  Update on gtest_child2 gtest_parent_2
  Update on gtest_child3 gtest_parent_3
-> Append  (cost=0.00..3.06 rows=3 width=14)
    -> Seq Scan on gtest_child gtest_parent_1  (cost=0.00..1.01 rows=1 width=14)
        Filter: (f2 = 101)
    -> Seq Scan on gtest_child2 gtest_parent_2  (cost=0.00..1.01 rows=1 width=14)
        Filter: (f2 = 101)
    -> Seq Scan on gtest_child3 gtest_parent_3  (cost=0.00..1.01 rows=1 width=14)
        Filter: (f2 = 101)
```

In this example the Update node needs to consider three child tables, but not the originally-mentioned partitioned table (since that never stores any data). So there are three input scanning subplans, one per table. For clarity, the Update node is annotated to show the specific target tables that will be updated, in the same order as the corresponding subplans.

The Planning time shown by EXPLAIN ANALYZE is the time it took to generate the query plan from the parsed query and optimize it. It does not include parsing or rewriting.

The Execution time shown by EXPLAIN ANALYZE includes executor start-up and shut-down time, as well as the time to run any triggers that are fired, but it does not include parsing, rewriting, or planning time. Time spent executing BEFORE triggers, if any, is included in the time for the related Insert, Update, or Delete node; but time spent executing AFTER triggers is not counted there because AFTER triggers are fired after completion of the whole plan. The total time spent in each trigger (either BEFORE or AFTER) is also shown separately. Note that deferred constraint triggers will not be executed until end of transaction and are thus not considered at all by EXPLAIN ANALYZE.

The time shown for the top-level node does not include any time needed to convert the query's output data into displayable form or to send it to the client. While EXPLAIN ANALYZE will never send the data to the client, it can be told to convert the query's output data to displayable form and measure the time needed for that, by specifying the SERIALIZE option. That time will be shown separately, and it's also included in the total Execution time.

14.1.3. Caveats

There are two significant ways in which run times measured by EXPLAIN ANALYZE can deviate from normal execution of the same query. First, since no output rows are delivered to the client, network transmission costs are not included. I/O conversion costs are not included either unless SERIALIZE is specified. Second, the measurement overhead added by EXPLAIN ANALYZE can be significant, especially on machines with slow `gettimeofday()` operating-system calls. You can use the [pg_test_timing](#) tool to measure the overhead of timing on your system.

EXPLAIN results should not be extrapolated to situations much different from the one you are actually testing; for example, results on a toy-sized table cannot be assumed to apply to large tables. The planner's cost estimates are not linear and so it might choose a different plan for a larger or smaller table. An extreme example is that on a table that only occupies one disk page, you'll nearly always get a sequential

scan plan whether indexes are available or not. The planner realizes that it's going to take one disk page read to process the table in any case, so there's no value in expending additional page reads to look at an index. (We saw this happening in the `polygon_tbl` example above.)

There are cases in which the actual and estimated values won't match up well, but nothing is really wrong. One such case occurs when plan node execution is stopped short by a `LIMIT` or similar effect. For example, in the `LIMIT` query we used before,

```
EXPLAIN ANALYZE SELECT * FROM tenk1 WHERE unique1 < 100 AND unique2 > 9000 LIMIT 2;
```

QUERY PLAN

```
-----
Limit  (cost=0.29..14.33 rows=2 width=244) (actual time=0.051..0.071 rows=2 loops=1)
->  Index Scan using tenk1_unique2 on tenk1  (cost=0.29..70.50 rows=10 width=244)
    (actual time=0.051..0.070 rows=2 loops=1)
        Index Cond: (unique2 > 9000)
        Filter: (unique1 < 100)
        Rows Removed by Filter: 287
Planning Time: 0.077 ms
Execution Time: 0.086 ms
```

the estimated cost and row count for the Index Scan node are shown as though it were run to completion. But in reality the Limit node stopped requesting rows after it got two, so the actual row count is only 2 and the run time is less than the cost estimate would suggest. This is not an estimation error, only a discrepancy in the way the estimates and true values are displayed.

Merge joins also have measurement artifacts that can confuse the unwary. A merge join will stop reading one input if it's exhausted the other input and the next key value in the one input is greater than the last key value of the other input; in such a case there can be no more matches and so no need to scan the rest of the first input. This results in not reading all of one child, with results like those mentioned for `LIMIT`. Also, if the outer (first) child contains rows with duplicate key values, the inner (second) child is backed up and rescanned for the portion of its rows matching that key value. `EXPLAIN ANALYZE` counts these repeated emissions of the same inner rows as if they were real additional rows. When there are many outer duplicates, the reported actual row count for the inner child plan node can be significantly larger than the number of rows that are actually in the inner relation.

BitmapAnd and BitmapOr nodes always report their actual row counts as zero, due to implementation limitations.

Normally, `EXPLAIN` will display every plan node created by the planner. However, there are cases where the executor can determine that certain nodes need not be executed because they cannot produce any rows, based on parameter values that were not available at planning time. (Currently this can only happen for child nodes of an Append or MergeAppend node that is scanning a partitioned table.) When this happens, those plan nodes are omitted from the `EXPLAIN` output and a `Subplans Removed: N` annotation appears instead.

14.2. Statistics Used by the Planner

14.2.1. Single-Column Statistics

As we saw in the previous section, the query planner needs to estimate the number of rows retrieved by a query in order to make good choices of query plans. This section provides a quick look at the statistics that the system uses for these estimates.

One component of the statistics is the total number of entries in each table and index, as well as the number of disk blocks occupied by each table and index. This information is kept in the table `pg_class`, in the columns `reltuples` and `relpages`. We can look at it with queries similar to this one:

```
SELECT relname, relkind, reltuples, relpages
```

```
FROM pg_class
WHERE relname LIKE 'tenk1%';
```

relname	relkind	reltuples	relpages
tenk1	r	10000	345
tenk1_hundred	i	10000	11
tenk1_thous_tenthous	i	10000	30
tenk1_unique1	i	10000	30
tenk1_unique2	i	10000	30

(5 rows)

Here we can see that `tenk1` contains 10000 rows, as do its indexes, but the indexes are (unsurprisingly) much smaller than the table.

For efficiency reasons, `reltuples` and `relpages` are not updated on-the-fly, and so they usually contain somewhat out-of-date values. They are updated by `VACUUM`, `ANALYZE`, and a few DDL commands such as `CREATE INDEX`. A `VACUUM` or `ANALYZE` operation that does not scan the entire table (which is commonly the case) will incrementally update the `reltuples` count on the basis of the part of the table it did scan, resulting in an approximate value. In any case, the planner will scale the values it finds in `pg_class` to match the current physical table size, thus obtaining a closer approximation.

Most queries retrieve only a fraction of the rows in a table, due to `WHERE` clauses that restrict the rows to be examined. The planner thus needs to make an estimate of the *selectivity* of `WHERE` clauses, that is, the fraction of rows that match each condition in the `WHERE` clause. The information used for this task is stored in the `pg_statistic` system catalog. Entries in `pg_statistic` are updated by the `ANALYZE` and `VACUUM ANALYZE` commands, and are always approximate even when freshly updated.

Rather than look at `pg_statistic` directly, it's better to look at its view `pg_stats` when examining the statistics manually. `pg_stats` is designed to be more easily readable. Furthermore, `pg_stats` is readable by all, whereas `pg_statistic` is only readable by a superuser. (This prevents unprivileged users from learning something about the contents of other people's tables from the statistics. The `pg_stats` view is restricted to show only rows about tables that the current user can read.) For example, we might do:

```
SELECT attname, inherited, n_distinct,
       array_to_string(most_common_vals, E'\n') as most_common_vals
FROM pg_stats
WHERE tablename = 'road';
```

attname	inherited	n_distinct	most_common_vals
name	f	-0.5681108	I- 580 Ramp+
			I- 880 Ramp+
			Sp Railroad +
			I- 580 +
			I- 680 Ramp+
			I- 80 Ramp+
			14th St +
			I- 880 +
			Mac Arthur Blvd+
			Mission Blvd+
...			
name	t	-0.5125	I- 580 Ramp+
			I- 880 Ramp+
			I- 580 +
			I- 680 Ramp+
			I- 80 Ramp+
			Sp Railroad +

```

      |           |           | I- 880           +
      |           |           | State Hwy 13     Ramp+
      |           |           | I- 80           +
      |           |           | State Hwy 24     Ramp+
...
thepath | f           |           | 0 |
thepath | t           |           | 0 |
(4 rows)

```

Note that two rows are displayed for the same column, one corresponding to the complete inheritance hierarchy starting at the `road` table (`inherited=t`), and another one including only the `road` table itself (`inherited=f`). (For brevity, we have only shown the first ten most-common values for the `name` column.)

The amount of information stored in `pg_statistic` by `ANALYZE`, in particular the maximum number of entries in the `most_common_vals` and `histogram_bounds` arrays for each column, can be set on a column-by-column basis using the `ALTER TABLE SET STATISTICS` command, or globally by setting the [default_statistics_target](#) configuration variable. The default limit is presently 100 entries. Raising the limit might allow more accurate planner estimates to be made, particularly for columns with irregular data distributions, at the price of consuming more space in `pg_statistic` and slightly more time to compute the estimates. Conversely, a lower limit might be sufficient for columns with simple data distributions.

Further details about the planner's use of statistics can be found in [Chapter 68](#).

14.2.2. Extended Statistics

It is common to see slow queries running bad execution plans because multiple columns used in the query clauses are correlated. The planner normally assumes that multiple conditions are independent of each other, an assumption that does not hold when column values are correlated. Regular statistics, because of their per-individual-column nature, cannot capture any knowledge about cross-column correlation. However, Postgres Pro has the ability to compute *multivariate statistics*, which can capture such information.

Because the number of possible column combinations is very large, it's impractical to compute multivariate statistics automatically. Instead, *extended statistics objects*, more often called just *statistics objects*, can be created to instruct the server to obtain statistics across interesting sets of columns.

Statistics objects are created using the `CREATE STATISTICS` command. Creation of such an object merely creates a catalog entry expressing interest in the statistics. Actual data collection is performed by `ANALYZE` (either a manual command, or background auto-analyze). The collected values can be examined in the `pg_statistic_ext_data` catalog.

`ANALYZE` computes extended statistics based on the same sample of table rows that it takes for computing regular single-column statistics. Since the sample size is increased by increasing the statistics target for the table or any of its columns (as described in the previous section), a larger statistics target will normally result in more accurate extended statistics, as well as more time spent calculating them.

The following subsections describe the kinds of extended statistics that are currently supported.

14.2.2.1. Functional Dependencies

The simplest kind of extended statistics tracks *functional dependencies*, a concept used in definitions of database normal forms. We say that column `b` is functionally dependent on column `a` if knowledge of the value of `a` is sufficient to determine the value of `b`, that is there are no two rows having the same value of `a` but different values of `b`. In a fully normalized database, functional dependencies should exist only on primary keys and superkeys. However, in practice many data sets are not fully normalized for various reasons; intentional denormalization for performance reasons is a common example. Even in a fully normalized database, there may be partial correlation between some columns, which can be expressed as partial functional dependency.

The existence of functional dependencies directly affects the accuracy of estimates in certain queries. If a query contains conditions on both the independent and the dependent column(s), the conditions on

the dependent columns do not further reduce the result size; but without knowledge of the functional dependency, the query planner will assume that the conditions are independent, resulting in underestimating the result size.

To inform the planner about functional dependencies, `ANALYZE` can collect measurements of cross-column dependency. Assessing the degree of dependency between all sets of columns would be prohibitively expensive, so data collection is limited to those groups of columns appearing together in a statistics object defined with the `dependencies` option. It is advisable to create `dependencies` statistics only for column groups that are strongly correlated, to avoid unnecessary overhead in both `ANALYZE` and later query planning.

Here is an example of collecting functional-dependency statistics:

```
CREATE STATISTICS stts (dependencies) ON city, zip FROM zipcodes;
```

```
ANALYZE zipcodes;
```

```
SELECT stxname, stxkeys, stxddependencies
FROM pg_statistic_ext join pg_statistic_ext_data on (oid = stxoid)
WHERE stxname = 'stts';
```

stxname	stxkeys	stxddependencies
stts	1 5	{"1 => 5": 1.000000, "5 => 1": 0.423130}

(1 row)

Here it can be seen that column 1 (zip code) fully determines column 5 (city) so the coefficient is 1.0, while city only determines zip code about 42% of the time, meaning that there are many cities (58%) that are represented by more than a single ZIP code.

When computing the selectivity for a query involving functionally dependent columns, the planner adjusts the per-condition selectivity estimates using the dependency coefficients so as not to produce an underestimate.

14.2.2.1.1. Limitations of Functional Dependencies

Functional dependencies are currently only applied when considering simple equality conditions that compare columns to constant values, and `IN` clauses with constant values. They are not used to improve estimates for equality conditions comparing two columns or comparing a column to an expression, nor for range clauses, `LIKE` or any other type of condition.

When estimating with functional dependencies, the planner assumes that conditions on the involved columns are compatible and hence redundant. If they are incompatible, the correct estimate would be zero rows, but that possibility is not considered. For example, given a query like

```
SELECT * FROM zipcodes WHERE city = 'San Francisco' AND zip = '94105';
```

the planner will disregard the `city` clause as not changing the selectivity, which is correct. However, it will make the same assumption about

```
SELECT * FROM zipcodes WHERE city = 'San Francisco' AND zip = '90210';
```

even though there will really be zero rows satisfying this query. Functional dependency statistics do not provide enough information to conclude that, however.

In many practical situations, this assumption is usually satisfied; for example, there might be a GUI in the application that only allows selecting compatible city and ZIP code values to use in a query. But if that's not the case, functional dependencies may not be a viable option.

14.2.2.2. Multivariate N-Distinct Counts

Single-column statistics store the number of distinct values in each column. Estimates of the number of distinct values when combining more than one column (for example, for `GROUP BY a, b`) are frequently wrong when the planner only has single-column statistical data, causing it to select bad plans.

To improve such estimates, `ANALYZE` can collect `n-distinct` statistics for groups of columns. As before, it's impractical to do this for every possible column grouping, so data is collected only for those groups of columns appearing together in a statistics object defined with the `ndistinct` option. Data will be collected for each possible combination of two or more columns from the set of listed columns.

Continuing the previous example, the `n-distinct` counts in a table of ZIP codes might look like the following:

```
CREATE STATISTICS stts2 (ndistinct) ON city, state, zip FROM zipcodes;

ANALYZE zipcodes;

SELECT stxkeys AS k, stxdndistinct AS nd
FROM pg_statistic_ext join pg_statistic_ext_data on (oid = stxoid)
WHERE stxname = 'stts2';
-[ RECORD 1 ]-----
k | 1 2 5
nd | {"1, 2": 33178, "1, 5": 33178, "2, 5": 27435, "1, 2, 5": 33178}
(1 row)
```

This indicates that there are three combinations of columns that have 33178 distinct values: ZIP code and state; ZIP code and city; and ZIP code, city and state (the fact that they are all equal is expected given that ZIP code alone is unique in this table). On the other hand, the combination of city and state has only 27435 distinct values.

It's advisable to create `ndistinct` statistics objects only on combinations of columns that are actually used for grouping, and for which misestimation of the number of groups is resulting in bad plans. Otherwise, the `ANALYZE` cycles are just wasted.

14.2.2.3. Multivariate MCV Lists

Another type of statistic stored for each column are most-common value lists. This allows very accurate estimates for individual columns, but may result in significant misestimates for queries with conditions on multiple columns.

To improve such estimates, `ANALYZE` can collect MCV lists on combinations of columns. Similarly to functional dependencies and `n-distinct` coefficients, it's impractical to do this for every possible column grouping. Even more so in this case, as the MCV list (unlike functional dependencies and `n-distinct` coefficients) does store the common column values. So data is collected only for those groups of columns appearing together in a statistics object defined with the `mcv` option.

Continuing the previous example, the MCV list for a table of ZIP codes might look like the following (unlike for simpler types of statistics, a function is required for inspection of MCV contents):

```
CREATE STATISTICS stts3 (mcv) ON city, state FROM zipcodes;

ANALYZE zipcodes;

SELECT m.* FROM pg_statistic_ext join pg_statistic_ext_data on (oid = stxoid),
pg_mcv_list_items(stxdmccv) m WHERE stxname = 'stts3';
```

index	values	nulls	frequency	base_frequency
0	{Washington, DC}	{f,f}	0.003467	2.7e-05
1	{Apo, AE}	{f,f}	0.003067	1.9e-05
2	{Houston, TX}	{f,f}	0.002167	0.000133
3	{El Paso, TX}	{f,f}	0.002	0.000113
4	{New York, NY}	{f,f}	0.001967	0.000114
5	{Atlanta, GA}	{f,f}	0.001633	3.3e-05
6	{Sacramento, CA}	{f,f}	0.001433	7.8e-05
7	{Miami, FL}	{f,f}	0.0014	6e-05

```

8 | {Dallas, TX}          | {f,f} | 0.001367 | 8.8e-05
9 | {Chicago, IL}         | {f,f} | 0.001333 | 5.1e-05
...
(99 rows)

```

This indicates that the most common combination of city and state is Washington in DC, with actual frequency (in the sample) about 0.35%. The base frequency of the combination (as computed from the simple per-column frequencies) is only 0.0027%, resulting in two orders of magnitude under-estimates.

It's advisable to create MCV statistics objects only on combinations of columns that are actually used in conditions together, and for which misestimation of the number of groups is resulting in bad plans. Otherwise, the `ANALYZE` and planning cycles are just wasted.

14.3. Controlling the Planner with Explicit JOIN Clauses

It is possible to control the query planner to some extent by using the explicit `JOIN` syntax. To see why this matters, we first need some background.

In a simple join query, such as:

```
SELECT * FROM a, b, c WHERE a.id = b.id AND b.ref = c.id;
```

the planner is free to join the given tables in any order. For example, it could generate a query plan that joins A to B, using the `WHERE` condition `a.id = b.id`, and then joins C to this joined table, using the other `WHERE` condition. Or it could join B to C and then join A to that result. Or it could join A to C and then join them with B — but that would be inefficient, since the full Cartesian product of A and C would have to be formed, there being no applicable condition in the `WHERE` clause to allow optimization of the join. (All joins in the Postgres Pro executor happen between two input tables, so it's necessary to build up the result in one or another of these fashions.) The important point is that these different join possibilities give semantically equivalent results but might have hugely different execution costs. Therefore, the planner will explore all of them to try to find the most efficient query plan.

When a query only involves two or three tables, there aren't many join orders to worry about. But the number of possible join orders grows exponentially as the number of tables expands. Beyond ten or so input tables it's no longer practical to do an exhaustive search of all the possibilities, and even for six or seven tables planning might take an annoyingly long time. When there are too many input tables, the Postgres Pro planner will switch from exhaustive search to a *genetic* probabilistic search through a limited number of possibilities. (The switch-over threshold is set by the `geqo_threshold` run-time parameter.) The genetic search takes less time, but it won't necessarily find the best possible plan.

When the query involves outer joins, the planner has less freedom than it does for plain (inner) joins. For example, consider:

```
SELECT * FROM a LEFT JOIN (b JOIN c ON (b.ref = c.id)) ON (a.id = b.id);
```

Although this query's restrictions are superficially similar to the previous example, the semantics are different because a row must be emitted for each row of A that has no matching row in the join of B and C. Therefore the planner has no choice of join order here: it must join B to C and then join A to that result. Accordingly, this query takes less time to plan than the previous query. In other cases, the planner might be able to determine that more than one join order is safe. For example, given:

```
SELECT * FROM a LEFT JOIN b ON (a.bid = b.id) LEFT JOIN c ON (a.cid = c.id);
```

it is valid to join A to either B or C first. Currently, only `FULL JOIN` completely constrains the join order. Most practical cases involving `LEFT JOIN` or `RIGHT JOIN` can be rearranged to some extent.

Explicit inner join syntax (`INNER JOIN`, `CROSS JOIN`, or unadorned `JOIN`) is semantically the same as listing the input relations in `FROM`, so it does not constrain the join order.

Even though most kinds of `JOIN` don't completely constrain the join order, it is possible to instruct the Postgres Pro query planner to treat all `JOIN` clauses as constraining the join order anyway. For example, these three queries are logically equivalent:

```
SELECT * FROM a, b, c WHERE a.id = b.id AND b.ref = c.id;
SELECT * FROM a CROSS JOIN b CROSS JOIN c WHERE a.id = b.id AND b.ref = c.id;
SELECT * FROM a JOIN (b JOIN c ON (b.ref = c.id)) ON (a.id = b.id);
```

But if we tell the planner to honor the `JOIN` order, the second and third take less time to plan than the first. This effect is not worth worrying about for only three tables, but it can be a lifesaver with many tables.

To force the planner to follow the join order laid out by explicit `JOINS`, set the `join_collapse_limit` run-time parameter to 1. (Other possible values are discussed below.)

You do not need to constrain the join order completely in order to cut search time, because it's OK to use `JOIN` operators within items of a plain `FROM` list. For example, consider:

```
SELECT * FROM a CROSS JOIN b, c, d, e WHERE ...;
```

With `join_collapse_limit = 1`, this forces the planner to join A to B before joining them to other tables, but doesn't constrain its choices otherwise. In this example, the number of possible join orders is reduced by a factor of 5.

Constraining the planner's search in this way is a useful technique both for reducing planning time and for directing the planner to a good query plan. If the planner chooses a bad join order by default, you can force it to choose a better order via `JOIN` syntax — assuming that you know of a better order, that is. Experimentation is recommended.

A closely related issue that affects planning time is collapsing of subqueries into their parent query. For example, consider:

```
SELECT *
FROM x, y,
    (SELECT * FROM a, b, c WHERE something) AS ss
WHERE somethingelse;
```

This situation might arise from use of a view that contains a join; the view's `SELECT` rule will be inserted in place of the view reference, yielding a query much like the above. Normally, the planner will try to collapse the subquery into the parent, yielding:

```
SELECT * FROM x, y, a, b, c WHERE something AND somethingelse;
```

This usually results in a better plan than planning the subquery separately. (For example, the outer `WHERE` conditions might be such that joining X to A first eliminates many rows of A, thus avoiding the need to form the full logical output of the subquery.) But at the same time, we have increased the planning time; here, we have a five-way join problem replacing two separate three-way join problems. Because of the exponential growth of the number of possibilities, this makes a big difference. The planner tries to avoid getting stuck in huge join search problems by not collapsing a subquery if more than `from_collapse_limit` `FROM` items would result in the parent query. You can trade off planning time against quality of plan by adjusting this run-time parameter up or down.

`from_collapse_limit` and `join_collapse_limit` are similarly named because they do almost the same thing: one controls when the planner will “flatten out” subqueries, and the other controls when it will flatten out explicit joins. Typically you would either set `join_collapse_limit` equal to `from_collapse_limit` (so that explicit joins and subqueries act similarly) or set `join_collapse_limit` to 1 (if you want to control join order with explicit joins). But you might set them differently if you are trying to fine-tune the trade-off between planning time and run time.

14.4. Populating a Database

One might need to insert a large amount of data when first populating a database. This section contains some suggestions on how to make this process as efficient as possible.

14.4.1. Disable Autocommit

When using multiple `INSERTs`, turn off autocommit and just do one commit at the end. (In plain SQL, this means issuing `BEGIN` at the start and `COMMIT` at the end. Some client libraries might do this behind your back, in which case you need to make sure the library does it when you want it done.) If you allow each insertion to be committed separately, Postgres Pro is doing a lot of work for each row that is added. An additional benefit of doing all insertions in one transaction is that if the insertion of one row were to fail then the insertion of all rows inserted up to that point would be rolled back, so you won't be stuck with partially loaded data.

14.4.2. Use `COPY`

Use `COPY` to load all the rows in one command, instead of using a series of `INSERT` commands. The `COPY` command is optimized for loading large numbers of rows; it is less flexible than `INSERT`, but incurs significantly less overhead for large data loads. Since `COPY` is a single command, there is no need to disable autocommit if you use this method to populate a table.

If you cannot use `COPY`, it might help to use `PREPARE` to create a prepared `INSERT` statement, and then use `EXECUTE` as many times as required. This avoids some of the overhead of repeatedly parsing and planning `INSERT`. Different interfaces provide this facility in different ways; look for “prepared statements” in the interface documentation.

Note that loading a large number of rows using `COPY` is almost always faster than using `INSERT`, even if `PREPARE` is used and multiple insertions are batched into a single transaction.

`COPY` is fastest when used within the same transaction as an earlier `CREATE TABLE` or `TRUNCATE` command. In such cases no WAL needs to be written, because in case of an error, the files containing the newly loaded data will be removed anyway. However, this consideration only applies when `wal_level` is `minimal` as all commands must write WAL otherwise.

14.4.3. Remove Indexes

If you are loading a freshly created table, the fastest method is to create the table, bulk load the table's data using `COPY`, then create any indexes needed for the table. Creating an index on pre-existing data is quicker than updating it incrementally as each row is loaded.

If you are adding large amounts of data to an existing table, it might be a win to drop the indexes, load the table, and then recreate the indexes. Of course, the database performance for other users might suffer during the time the indexes are missing. One should also think twice before dropping a unique index, since the error checking afforded by the unique constraint will be lost while the index is missing.

14.4.4. Remove Foreign Key Constraints

Just as with indexes, a foreign key constraint can be checked “in bulk” more efficiently than row-by-row. So it might be useful to drop foreign key constraints, load data, and re-create the constraints. Again, there is a trade-off between data load speed and loss of error checking while the constraint is missing.

What's more, when you load data into a table with existing foreign key constraints, each new row requires an entry in the server's list of pending trigger events (since it is the firing of a trigger that checks the row's foreign key constraint). Loading many millions of rows can cause the trigger event queue to overflow available memory, leading to intolerable swapping or even outright failure of the command. Therefore it may be *necessary*, not just desirable, to drop and re-apply foreign keys when loading large amounts of data. If temporarily removing the constraint isn't acceptable, the only other recourse may be to split up the load operation into smaller transactions.

14.4.5. Increase `maintenance_work_mem`

Temporarily increasing the `maintenance_work_mem` configuration variable when loading large amounts of data can lead to improved performance. This will help to speed up `CREATE INDEX` commands and `ALTER TABLE ADD FOREIGN KEY` commands. It won't do much for `COPY` itself, so this advice is only useful when you are using one or both of the above techniques.

14.4.6. Increase `max_wal_size`

Temporarily increasing the `max_wal_size` configuration variable can also make large data loads faster. This is because loading a large amount of data into Postgres Pro will cause checkpoints to occur more often than the normal checkpoint frequency (specified by the `checkpoint_timeout` configuration variable). Whenever a checkpoint occurs, all dirty pages must be flushed to disk. By increasing `max_wal_size` temporarily during bulk data loads, the number of checkpoints that are required can be reduced.

14.4.7. Run `shardman.analyze()` Afterwards

Whenever you have significantly altered the distribution of data within a table, running `shardman.analyze()` is strongly recommended. This includes bulk loading large amounts of data into the table. Running `shardman.analyze()` (or `VACUUM ANALYZE`) ensures that the planner has up-to-date statistics about the table. With no statistics or obsolete statistics, the planner might make poor decisions during query planning, leading to poor performance on any tables with inaccurate or nonexistent statistics. Note that if the autovacuum daemon is enabled, it might run `shardman.analyze()` automatically; see [Section 24.1.3](#) and [Section 24.1.6](#) for more information.

14.4.8. Some Notes about `pg_dump`

Dump scripts generated by `pg_dump` automatically apply several, but not all, of the above guidelines. To restore a `pg_dump` dump as quickly as possible, you need to do a few extra things manually. (Note that these points apply while *restoring* a dump, not while *creating* it. The same points apply whether loading a text dump with `psql` or using `pg_restore` to load from a `pg_dump` archive file.)

By default, `pg_dump` uses `COPY`, and when it is generating a complete schema-and-data dump, it is careful to load data before creating indexes and foreign keys. So in this case several guidelines are handled automatically. What is left for you to do is to:

- Set appropriate (i.e., larger than normal) values for `maintenance_work_mem` and `max_wal_size`.
- If using WAL archiving or streaming replication, consider disabling them during the restore. To do that, set `archive_mode` to `off`, `wal_level` to `minimal`, and `max_wal_senders` to zero before loading the dump. Afterwards, set them back to the right values and take a fresh base backup.
- Experiment with the parallel dump and restore modes of both `pg_dump` and `pg_restore` and find the optimal number of concurrent jobs to use. Dumping and restoring in parallel by means of the `-j` option should give you a significantly higher performance over the serial mode.
- Consider whether the whole dump should be restored as a single transaction. To do that, pass the `-1` or `--single-transaction` command-line option to `psql` or `pg_restore`. When using this mode, even the smallest of errors will rollback the entire restore, possibly discarding many hours of processing. Depending on how interrelated the data is, that might seem preferable to manual cleanup, or not. `COPY` commands will run fastest if you use a single transaction and have WAL archiving turned off.
- If multiple CPUs are available in the database server, consider using `pg_restore`'s `--jobs` option. This allows concurrent data loading and index creation.
- Run `ANALYZE` afterwards.

A data-only dump will still use `COPY`, but it does not drop or recreate indexes, and it does not normally touch foreign keys.¹ So when loading a data-only dump, it is up to you to drop and recreate indexes and foreign keys if you wish to use those techniques. It's still useful to increase `max_wal_size` while loading the data, but don't bother increasing `maintenance_work_mem`; rather, you'd do that while manually recreating indexes and foreign keys afterwards. And don't forget to `ANALYZE` when you're done; see [Section 24.1.3](#) and [Section 24.1.6](#) for more information.

14.5. Autoprepared Statements

¹ You can get the effect of disabling foreign keys by using the `--disable-triggers` option — but realize that that eliminates, rather than just postpones, foreign key validation, and so it is possible to insert bad data if you use it.

Postgres Pro Shardman provides the *autoprepare* mode that can implicitly prepare frequently used statements to eliminate the cost of their compilation and planning on each subsequent execution.

While providing slightly lower performance gains than explicitly prepared statements, this mode enables you to get prepared plans for cases where explicit `PREPARE` commands are not supported, such as using `pgbouncer` with a pooling level other than `session`, working with distributed partitioned tables, or running database applications that are not designed to execute prepared statements.

By default, the `autoprepare` mode is switched off. To enable this feature, assign a non-zero value to the `autoprepare_threshold` configuration variable, which sets the minimal number of times a statement should be executed to get autoprepared. Once the same query with different literal values is repeated the specified number of times, Postgres Pro builds a generic plan for this query, replacing all constant literals with parameters. The generic plan will be chosen for execution if the estimated time of this plan is lower than an average time of the customized plans. Just like for explicitly prepared statements, generic plans get invalidated when a catalog change occurs.

When this feature is enabled, queries submitted via both simple and extended query protocols could be autoprepared. The extended protocol supports the transmission of parameterized queries, and autopreparing of such queries incurs less overhead. On the other hand, queries received via different protocols can be very diverse in nature, and autopreparing is ineffective for some query types. Hence, the `autoprepare_for_protocol` parameter allows to enable/disable the statement autopreparing on the protocol level.

If your application issues many different statements, autopreparing them all can cause memory overflow unless you impose any cache restrictions. It is especially important when running multiple active clients as the cache of autoprepared statements is local to the backend. To avoid cache bloat, you can do the following:

- Limit the number of autoprepared statements per backend using the `autoprepare_limit` parameter. This setting is recommended for workloads with multiple simple queries.
- Limit the amount of memory that can be allocated for autoprepared statements on a backend using the `autoprepare_memory_limit` parameter. This setting can incur some overhead when calculating the amount of cache used for autoprepared statements, so it is only recommended for workloads that contain complex queries for which limiting the number of autoprepared statements may be ineffective.

In both cases, most frequently used queries are kept in memory using the LRU strategy. If both parameters are set, Postgres Pro enforces the first reached limit.

Note that the cache of autoprepared statements is fully cleaned up when executing the following commands:

- `DISCARD PLANS`
- `DISCARD ALL`
- Any `ALTER SYSTEM...` command
- Any `SET var` command, both standalone and inside other commands. For example:

```
ALTER TABLE test ALTER COLUMN a SET STORAGE PLAIN
```
- Any DDL command that supports event triggers, such as `CREATE/DROP TABLE`

You can check all autoprepared statements available in the current session in the `pg_autoprepared_statements` system view. It shows the original text of the query, types of the extracted parameters that replace literals, and the query execution counter.

14.6. Postgres Pro Shardman Performance Tuning

Below is a list of performance tuning tips for the distributed systems.

14.6.1. Examining Plans

Tuning query execution is better on a subset of production data that represents actual data distribution. Let's look at some sample plans.

```
EXPLAIN VERBOSE
SELECT bid, avg(abalance) FROM pgbench_accounts
WHERE bid IN (10,20,30,40)
GROUP BY bid;
```

QUERY PLAN

```

Append (cost=0.29..21.98 rows=4 width=36)
-> GroupAggregate (cost=0.29..18.98 rows=1 width=36)
    Output: pgbench_accounts.bid, avg(pgbench_accounts.abalance)
    Group Key: pgbench_accounts.bid
-> Index Scan using pgbench_accounts_15_pkey on public.pgbench_accounts_15
pgbench_accounts (cost=0.29..18.96 rows=1 width=8)
    Output: pgbench_accounts.bid, pgbench_accounts.abalance
    Index Cond: (pgbench_accounts.bid = ANY ('{10,20,30,40}'::integer[]))
-> Async Foreign Scan (cost=0.99..0.99 rows=1 width=36)
    Output: pgbench_accounts_1.bid, (avg(pgbench_accounts_1.abalance))
    Relations: Aggregate on (public.pgbench_accounts_16_fdw pgbench_accounts_1)
    Remote SQL: SELECT bid, avg(abalance) FROM public.pgbench_accounts_16 WHERE
((bid = ANY ('{10,20,30,40}'::integer[]))) GROUP BY 1
    Transport: Silk
-> Async Foreign Scan (cost=0.99..0.99 rows=1 width=36)
    Output: pgbench_accounts_2.bid, (avg(pgbench_accounts_2.abalance))
    Relations: Aggregate on (public.pgbench_accounts_17_fdw pgbench_accounts_2)
    Remote SQL: SELECT bid, avg(abalance) FROM public.pgbench_accounts_17 WHERE
((bid = ANY ('{10,20,30,40}'::integer[]))) GROUP BY 1
    Transport: Silk
-> Async Foreign Scan (cost=1.00..1.00 rows=1 width=36)
    Output: pgbench_accounts_3.bid, (avg(pgbench_accounts_3.abalance))
    Relations: Aggregate on (public.pgbench_accounts_19_fdw pgbench_accounts_3)
    Remote SQL: SELECT bid, avg(abalance) FROM public.pgbench_accounts_19 WHERE
((bid = ANY ('{10,20,30,40}'::integer[]))) GROUP BY 1
    Transport: Silk
Query Identifier: -1714706980364121548

```

We see here that queries scanning three partitions are going to be sent to other nodes, coordinator data is also going to be scanned using `Index Scan`. We do not know what plan will be used on the remote side, but we see which queries will be sent (marked with `Remote SQL`). Note that `Transport: Silk` section is present in the foreign scan description. This indicates that Silk transport will be used to transfer results. We see that `Async` foreign scan is going to be used, which is fine. To discover which servers are used in the query, we should look at foreign tables definitions. For example, we can find out that `public.pgbench_accounts_19_fdw` is located on the `shardman_rq_2` server listening on `127.0.0.2:65432`:

```
SELECT srvname,srvoptions FROM pg_foreign_server s JOIN pg_foreign_table ON ftserver =
    s.oid
WHERE ftrelid = 'public.pgbench_accounts_19_fdw'::regclass;
-[ RECORD
   1 ]-----
srvname      | shardman_rg_2
srvoptions   | {async_capable=on,batch_size=100,binary_format=on,connect_timeout=5,dbname=postgres,extended_gin_index=off}
```

Now we can connect to `shardman_rg_2` server and find out which plan is used for the local query which was shown by the above `EXPLAIN`:

```
EXPLAIN SELECT bid, avg(abalance)
```

```
FROM public.pgbench_accounts_19
WHERE ((bid = ANY ('{10,20,30,40}'::integer[]))) GROUP BY 1;
```

QUERY PLAN

```
HashAggregate (cost=3641.00..3641.01 rows=1 width=36)
  Group Key: bid
    -> Seq Scan on pgbench_accounts_19 (cost=0.00..3141.00 rows=100000 width=8)
        Filter: (bid = ANY ('{10,20,30,40}'::integer[]))
```

While looking at distributed query plans, we can see that sometimes aggregates are not pushed down:

```
EXPLAIN VERBOSE
SELECT avg(abalance) FROM pgbench_accounts;
```

QUERY PLAN

```
Finalize Aggregate (cost=156209.38..156209.39 rows=1 width=32) (actual
time=590.359..590.371 rows=1 loops=1)
  Output: avg(pgbench_accounts.abalance)
    -> Append (cost=2891.00..156209.33 rows=20 width=32) (actual time=56.815..590.341
rows=20 loops=1)
      -> Partial Aggregate (cost=2891.00..2891.01 rows=1 width=32) (actual
time=56.812..56.813 rows=1 loops=1)
        Output: PARTIAL avg(pgbench_accounts.abalance)
          -> Seq Scan on public.pgbench_accounts_0 pgbench_accounts
(cost=0.00..2641.00 rows=100000 width=4) (actual time=0.018..38.478 rows=100000
loops=1)
            Output: pgbench_accounts.abalance
          -> Partial Aggregate (cost=23991.00..23991.01 rows=1 width=32) (actual
time=75.133..75.134 rows=1 loops=1)
            Output: PARTIAL avg(pgbench_accounts_1.abalance)
              -> Foreign Scan on public.pgbench_accounts_1_fdw pgbench_accounts_1
(cost=100.00..23741.00 rows=100000 width=4) (actual time=41.281..67.293 rows=100000
loops=1)
                Output: pgbench_accounts_1.abalance
                Remote SQL: SELECT abalance FROM public.pgbench_accounts_1
                Transport: Silk
.....
```

Here `avg()` is calculated on the coordinator side. This can lead to a significant growth of data transfer between nodes. The actual data transfer can be monitored with the `NETWORK` parameter of `EXPLAIN ANALYZE` (look at the `Network received` field of the topmost plan node):

```
EXPLAIN (ANALYZE, VERBOSE, NETWORK)
SELECT avg(abalance) FROM pgbench_accounts
```

QUERY PLAN

```
Finalize Aggregate (cost=156209.38..156209.39 rows=1 width=32) (actual
time=589.014..589.027 rows=1 loops=1)
  Output: avg(pgbench_accounts.abalance)
  Network: FDW bytes sent=3218 received=14402396
    -> Append (cost=2891.00..156209.33 rows=20 width=32) (actual time=52.111..588.999
rows=20 loops=1)
      Network: FDW bytes sent=3218 received=14402396
        -> Partial Aggregate (cost=2891.00..2891.01 rows=1 width=32) (actual
time=52.109..52.109 rows=1 loops=1)
```

```
Output: PARTIAL avg(pgbench_accounts.abalance)
-> Seq Scan on public.pgbench_accounts_0 pgbench_accounts
(cost=0.00..2641.00 rows=100000 width=4) (actual time=0.020..34.472 rows=100000
loops=1)

Output: pgbench_accounts.abalance
-> Partial Aggregate (cost=23991.00..23991.01 rows=1 width=32) (actual
time=78.616..78.617 rows=1 loops=1)
Output: PARTIAL avg(pgbench_accounts_1.abalance)
Network: FDW bytes sent=247 received=2400360
-> Foreign Scan on public.pgbench_accounts_1_fdw pgbench_accounts_1
(cost=100.00..23741.00 rows=100000 width=4) (actual time=42.359..69.984 rows=100000
loops=1)

Output: pgbench_accounts_1.abalance
Remote SQL: SELECT abalance FROM public.pgbench_accounts_1
Transport: Silk
Network: FDW bytes sent=247 received=2400360

.....
```

In such cases, we sometimes can rewrite the query:

```
EXPLAIN (ANALYZE, NETWORK, VERBOSE)
SELECT sum(abalance)::float/count(abalance) FROM pgbench_accounts where abalance is not
null;
```

QUERY PLAN

```
Finalize Aggregate (cost=12577.20..12577.22 rows=1 width=8) (actual
time=151.632..151.639 rows=1 loops=1)
Output: ((sum(pgbench_accounts.abalance))::double precision /
(count(pgbench_accounts.abalance))::double precision)
Network: FDW bytes sent=3907 received=872
-> Append (cost=3141.00..12577.10 rows=20 width=16) (actual time=55.589..151.621
rows=20 loops=1)
Network: FDW bytes sent=3907 received=872
-> Partial Aggregate (cost=3141.00..3141.01 rows=1 width=16) (actual
time=55.423..55.424 rows=1 loops=1)
Output: PARTIAL sum(pgbench_accounts.abalance), PARTIAL
count(pgbench_accounts.abalance)
-> Seq Scan on public.pgbench_accounts_0 pgbench_accounts
(cost=0.00..2641.00 rows=100000 width=4) (actual time=0.023..37.212 rows=100000
loops=1)

Output: pgbench_accounts.abalance
Filter: (pgbench_accounts.abalance IS NOT NULL)
-> Async Foreign Scan (cost=1.00..1.00 rows=1 width=16) (actual
time=0.055..0.089 rows=1 loops=1)
Output: (PARTIAL sum(pgbench_accounts_1.abalance)), (PARTIAL
count(pgbench_accounts_1.abalance))
Relations: Aggregate on (public.pgbench_accounts_1_fdw
pgbench_accounts_1)
Remote SQL: SELECT sum(abalance), count(abalance) FROM
public.pgbench_accounts_1 WHERE ((abalance IS NOT NULL))
Transport: Silk
Network: FDW bytes sent=300 received=800

.....
```

Rewriting the query here, we could decrease incoming network traffic generated by the query from 13 MB to 872 bytes.

Now let's look at two nearly identical joins.

```
EXPLAIN ANALYZE SELECT count(*) FROM pgbench_branches b
JOIN pgbench_history h ON b.bid = h.bid
WHERE mtime > '2023-03-14 10:00:00'::timestampz AND b.bbalance > 0;
```

QUERY PLAN

```
-----
Finalize Aggregate  (cost=8125.68..8125.69 rows=1 width=8) (actual time=27.464..27.543
rows=1 loops=1)
  -> Append  (cost=3.85..8125.63 rows=20 width=8) (actual time=0.036..27.475 rows=20
loops=1)
    -> Partial Aggregate  (cost=3.85..3.86 rows=1 width=8) (actual
time=0.033..0.036 rows=1 loops=1)
      -> Nested Loop  (cost=0.00..3.69 rows=67 width=0) (actual
time=0.025..0.027 rows=0 loops=1)
        Join Filter: (b.bid = h.bid)
        -> Seq Scan on pgbench_branches_0 b  (cost=0.00..1.01 rows=1
width=4) (actual time=0.023..0.024 rows=0 loops=1)
          Filter: (bbalance > 0)
          Rows Removed by Filter: 1
        -> Seq Scan on pgbench_history_0 h  (cost=0.00..1.84 rows=67
width=4) (never executed)
          Filter: (mtime > '2023-03-14 10:00:00+03'::timestamp with
time zone)
      -> Partial Aggregate  (cost=222.65..222.66 rows=1 width=8) (actual
time=3.969..3.973 rows=1 loops=1)
        -> Nested Loop  (cost=200.00..222.43 rows=86 width=0) (actual
time=3.736..3.920 rows=86 loops=1)
          Join Filter: (b_1.bid = h_1.bid)
          -> Foreign Scan on pgbench_branches_1_fdw b_1
(cost=100.00..101.22 rows=1 width=4) (actual time=1.929..1.932 rows=1 loops=1)
          -> Foreign Scan on pgbench_history_1_fdw h_1
(cost=100.00..120.14 rows=86 width=4) (actual time=1.795..1.916 rows=86 loops=1)
            Filter: (mtime > '2023-03-14 10:00:00+03'::timestamp with
time zone)
        -> Partial Aggregate  (cost=864.54..864.55 rows=1 width=8) (actual
time=1.780..1.786 rows=1 loops=1)
          -> Hash Join  (cost=200.01..864.53 rows=5 width=0) (actual
time=1.769..1.773 rows=0 loops=1)
            Hash Cond: (h_2.bid = b_2.bid)
            -> Foreign Scan on pgbench_history_2_fdw h_2
(cost=100.00..760.81 rows=975 width=4) (never executed)
              Filter: (mtime > '2023-03-14 10:00:00+03'::timestamp with
time zone)
            -> Hash  (cost=100.00..100.00 rows=1 width=4) (actual
time=1.740..1.742 rows=0 loops=1)
              Buckets: 1024  Batches: 1  Memory Usage: 8kB
              -> Foreign Scan on pgbench_branches_2_fdw b_2
(cost=100.00..100.00 rows=1 width=4) (actual time=1.738..1.738 rows=0 loops=1)
      ....
Planning Time: 6.066 ms
Execution Time: 33.851 ms
```

An interesting thing to note is that joining of `pgbench_branches` and `pgbench_history` partitions happens locally. It is a fetch-all plan — you can discover this by joins being located above foreign scans. It is not always evident why join pushdown does not happen. But if we look at the `pgbench_history` definition, we can see that `mtime` has the timestamp without time zone type.

```
\d pgbench_history
```

```

Partitioned table "public.pgbench_history"
Column |          Type          | Collation | Nullable | Default
-----+-----+-----+-----+-----
tid    | integer                |           |          |
bid    | integer                |           |          |
aid    | integer                |           |          |
delta  | integer                |           |          |
mtime  | timestamp without time zone |           |          |
filler | character(22)           |           |          |
Partition key: HASH (bid)
Number of partitions: 20 (Use \d+ to list them.)

```

And in the above query, the string describing time is converted to `timestamp with timezone`. This requires comparison of `mtime` column (of `timestamp` type) and `timestampz` value. The comparison is implicitly performed using the stable function `timestamp_gt_timestampz`. A filter containing a non-immutable function cannot be pushed down to the foreign server, so join is executed locally. If we rewrite the query, converting the string to a timestamp, we can see not only that joins are pushed down, but also that remote queries can be executed asynchronously because foreign scans in a plan tree are located immediately below Append:

```

EXPLAIN ANALYZE SELECT count(*) FROM pgbench_branches b
JOIN pgbench_history h ON b.bid = h.bid
WHERE mtime > '2023-03-14 10:00:00'::timestamp AND b.bbalance > 0;

```

QUERY PLAN

```

Finalize Aggregate  (cost=84.30..84.31 rows=1 width=8) (actual time=22.962..22.990
rows=1 loops=1)
  -> Append  (cost=3.85..84.25 rows=20 width=8) (actual time=0.196..22.927 rows=20
loops=1)
    -> Partial Aggregate  (cost=3.85..3.86 rows=1 width=8) (actual
time=0.032..0.034 rows=1 loops=1)
      -> Nested Loop  (cost=0.00..3.69 rows=67 width=0) (actual
time=0.024..0.026 rows=0 loops=1)
        Join Filter: (b.bid = h.bid)
        -> Seq Scan on pgbench_branches_0 b  (cost=0.00..1.01 rows=1
width=4) (actual time=0.023..0.023 rows=0 loops=1)
          Filter: (bbalance > 0)
          Rows Removed by Filter: 1
        -> Seq Scan on pgbench_history_0 h  (cost=0.00..1.84 rows=67
width=4) (never executed)
          Filter: (mtime > '2023-03-14 10:00:00'::timestamp without
time zone)
            -> Async Foreign Scan  (cost=0.99..0.99 rows=1 width=8) (actual
time=10.870..10.871 rows=1 loops=1)
              Relations: Aggregate on ((pgbench_branches_1_fdw b_1) INNER JOIN
(pgbench_history_1_fdw h_1))
            -> Async Foreign Scan  (cost=0.99..0.99 rows=1 width=8) (actual
time=0.016..0.017 rows=1 loops=1)
              Relations: Aggregate on ((pgbench_branches_2_fdw b_2) INNER JOIN
(pgbench_history_2_fdw h_2))
...
Planning Time: 7.729 ms
Execution Time: 14.603 ms

```

Note that foreign scans here include a list of joined relations. The expected cost of a foreign join is below 1.0. This is due to an optimistic technique of foreign join cost estimation, turned on by the `postgres_fdw.enforce_foreign_join` setting. Compare the total execution time (planning time + execution time) of the original and modified query — we could decrease it from about 40 to 22 ms.

Overall, while examining query plans, pay attention to what queries are actually pushed down. Some of the common reasons why joins cannot be pushed down is the absence of equi-joins on the sharding key and filters that contain non-immutable functions (possibly implicitly). If data is fetched from multiple replication groups, check that execution is mostly asynchronous.

14.6.1.1. EXPLAIN Parameters

This section lists `EXPLAIN` parameters for a distributed system.

`NETWORK` (boolean)

Include the actual data transfer between nodes in the `EXPLAIN ANALYZE` output. If this parameter is not specified, `off` is assumed. If the parameter is specified without a value, `on` is assumed.

`REMOTE` (boolean)

Include plans for queries executed on foreign servers. If this parameter or its value is not specified, `on` is assumed.

14.6.2. DML Optimizations

While evaluating performance of DML statements, it is important to understand how they are processed in Postgres Pro Shardman.

First of all, the execution of `INSERT` significantly differs from the execution of `UPDATE` and `DELETE` statements. The behavior of `INSERT` for sharded tables is controlled by the `batch_size` foreign server option, which can be set in `FDWOptions` section of [Postgres Pro Shardman configuration file](#). If `batch_size` is greater than 0, an `INSERT` in the same statement of several values that fall into the same foreign partition leads to the values being grouped together in batches of the specified size. Remote `INSERT` statements are prepared with the necessary number of parameters and then are executed with the given values. If the number of values does not match the number of prepared arguments, the modified statement with the necessary number of parameters is prepared again. A batch insert optimization can fail if a transaction inserts records one by one or records routed to different foreign tables are intermixed in one `INSERT` statement. A batch is formed for a single foreign modify operation. It is sent to the remote server when the batch is filled or when the modify operation is over. The modify operation is over when we start routing tuples to another sharded table partition. So, for bulk load, inserting multiple values in a single `INSERT` command or using `COPY` is recommended (as `COPY` is optimized in a similar way). Large `batch_size` values allow issuing less `INSERT` statements on remote side and so significantly reduce communication costs. However, during construction of parameters for prepared `INSERT` statements, all inserted values should be copied to libpq-allocated memory. This can lead to unrestricted memory usage on the query coordinator side when several large `text` or `bytea` objects are loaded.

`UPDATE` and `DELETE` statements can be executed in a direct or indirect mode. A direct mode is used when a statement can be directly sent to a foreign server. In this mode, to modify a table on a remote server, a new statement is created based on the original `ModifyTable` plan node. Using a direct update is not always possible. In particular, it is impossible when some conditions should be evaluated locally. In this case, a much less efficient indirect modification is used. An indirect modification includes several statements. The first one is `SELECT FOR UPDATE` to lock remote rows. The second one is an actual `UPDATE` or `DELETE`, which is prepared once and then executed with different parameters for each row of the `SELECT FOR UPDATE` statement result after local filters are applied to the result. Evidently, direct modifications are much more efficient.

You can easily identify whether a DML statement is going to be executed in a direct or indirect mode looking at the query plan. A typical example of an indirect modification is:

```
EXPLAIN VERBOSE DELETE FROM pgbench_history
WHERE bid = 20 AND mtime > '2023-03-14 10:00:00'::timestampz;
                                QUERY PLAN
```

```
Delete on public.pgbench_history (cost=100.00..142.66 rows=0 width=0)
  Foreign Delete on public.pgbench_history_17_fdw pgbench_history_1
    Remote SQL: DELETE FROM public.pgbench_history_17 WHERE ctid = $1
-> Foreign Scan on public.pgbench_history_17_fdw pgbench_history_1
(cost=100.00..142.66 rows=4 width=10)
  Output: pgbench_history_1.tableoid, pgbench_history_1.ctid
  Filter: (pgbench_history_1.mtime > '2023-03-14 10:00:00+03'::timestamp with
time zone)
    Remote SQL: SELECT mtime, ctid FROM public.pgbench_history_17 WHERE ((bid =
20)) FOR UPDATE
```

If we had chosen another type for the string constant, this would become a direct update.

```
EXPLAIN VERBOSE DELETE FROM pgbench_history
WHERE bid = 20 AND mtime > '2023-03-14 10:00:00'::timestamp;
explain verbose delete from pgbench_history where bid = 20 and mtime > '2023-03-14
10:00:00'::timestamp;
```

QUERY PLAN

```
Delete on public.pgbench_history (cost=100.00..146.97 rows=0 width=0)
  Foreign Delete on public.pgbench_history_17_fdw pgbench_history_1
-> Foreign Delete on public.pgbench_history_17_fdw pgbench_history_1
(cost=100.00..146.97 rows=4 width=10)
  Remote SQL: DELETE FROM public.pgbench_history_17 WHERE ((mtime > '2023-03-14
10:00:00'::timestamp without time zone)) AND ((bid = 20))
```

We see that in a direct update mode, only one statement is executed on the remote server.

14.6.2.1. DML Optimizations of Global Tables

The [shardman.gt_batch_size](#) configuration parameter, which you can tune, defines the size of an intermediate buffer used before sending data to a remote server.

INSERT uses the binary protocol and creates batches of the `shardman.gt_batch_size` size. Large values of the buffer size enable sending fewer network requests on the remote side and thus substantially reduce the connection costs. On the other hand, large values of this parameter can increase memory consumption on the query coordinator side. Therefore, when specifying the buffer size, it is important to achieve a compromise between the connection costs and the allocated memory size.

For UPDATE, a query for each column and each row is created on the coordinator and sent to remote nodes.

For DELETE, a query for a batch of data of the `shardman.gt_batch_size` size is created on the coordinator and sent to remote nodes.

14.6.3. Time Synchronization

The algorithm that provides data consistency on all the cluster nodes uses the system clock installed on the hosts. Therefore, the transaction commit latency depends on clock drift on different hosts, as the coordinator always waits for the most lagging host to catch up. This makes it crucial that the time on all the connected nodes of a Postgres Pro Shardman cluster are synchronized, as lack of synchronization may have a negative impact on Postgres Pro Shardman performance by increasing the query latency.

First, to ensure time synchronization on all cluster nodes, install chrony daemon when deploying a new cluster.

```
sudo apt update
sudo apt install -y chrony
sudo systemctl enable --now chrony
```

Check that chrony is working properly.

chronyc tracking

Expected output:

```
Reference ID      : C0248F82 (Time100.Stupi.SE)
Stratum          : 2
Ref time (UTC)   : Tue Apr 18 11:50:44 2023
System time      : 0.000019457 seconds slow of NTP time
Last offset      : -0.000005579 seconds
RMS offset       : 0.000089375 seconds
Frequency        : 30.777 ppm fast
Residual freq    : -0.000 ppm
Skew             : 0.003 ppm
Root delay       : 0.018349268 seconds
Root dispersion  : 0.000334640 seconds
Update interval  : 1039.1 seconds
Leap status      : Normal
```

Note that managing the clock drift should be performed using the OS tools. Postgres Pro Shardman diagnostic tools cannot be considered as the only and defining measurement utility.

To see if any major drift already exists, use the [shardman.pg_stat_csn](#) view that shows statistics on delays that take place during import of CSN snapshots. Its values are calculated when any related action is performed, or if any of the `shardman.trim_csnxid_map()` or `shardman.pg_oldest_csn_snapshot()` functions are called. These functions are called from the `csn trimmer routine worker`, therefore disabling this worker will result in these statistics not being collected.

The `csn_max_shift` field of the `shardman.pg_stat_csn` view shows the maximum registered snapshot CSN shift that caused a delay. This value defines the clock drift between the nodes in the cluster. A consecutive increase of this value means at least one's cluster system clock is out of sync. If this value exceeds 1000 (microseconds), it is recommended to check the time synchronization settings.

The same can be discovered if the `csn_total_import_delay` value increases while `csn_max_shift` remains unchanged. However, one-time increase may be due to single failures, non-related to the time issues.

Also, if the difference between `CSNXidMap_head_csn` and `shardman.oldest_csn` exceeds the `csn_snapshot_defer_time` parameter value and stays the same for a long time, it means that the `CSNSnapshotXidMap` map is full. It can result in a global transaction failure.

There are two main reasons for this issue.

- There is a transaction that runs for more than `csn_snapshot_defer_time` seconds and holds the entire cluster, holding the `VACUUM` process. In this case, `xid` field of the `shardman.oldest_csn` view is used to determine the transaction ID of this transaction, and the `rgid` field is used to determine the cluster node where this transaction is located.
- The `CSNSnapshotXidMap` map lacks capacity. During the normal operation the system might have transactions that exceed the `csn_snapshot_defer_time` value. To fix it, increase the `csn_snapshot_defer_time` time so that these transactions stay below this value.

If the `shardman.silk_tracepoints` configuration parameter is enabled, executing the `EXPLAIN` command for the distributed queries outputs the rows with information about how much time was spent on the query execution and what result it ended with, depending on the system components. These rows show metric values for the time spent on each component. The `net (qry)`, `net (1st tup)`, `net (last tup)` metrics calculate the difference between timestamps on different servers. This difference includes both time spent on a message transfer and the clock drift (positive or negative) between these servers. Therefore, these metrics can also help to determine whether there is any clock drift.

14.6.4. Distributed Query Diagnostics

Postgres Pro Shardman enhances the `EXPLAIN` command so that it can provide additional information about a query if it is distributed. The work with the distributed tables is based on the plan nodes with the `ForeignScan` type. A query to each remote partition is determined by a single plan node of this type, with Postgres Pro Shardman submitting additional information to the `EXPLAIN` blocks with the node description.

When executing a distributed query, the part of the plan (a subtree) that relates to a specific remote partition is serialized into an SQL statement. This process is known as deparsing. Then, this statement is sent to a remote server. The result of this query is the output of a `ForeignScan` node. It is used to gather the final results of the distributed query execution.

When the `VERBOSE` option of the `EXPLAIN` command is set to `on`, the `Remote SQL` field of the `ForeignScan` node block shows the statement sent to the remote server. Also, the `Server` field indicates the name of the server as it was specified during the cluster configuration and as it is displayed in `pg_foreign_server`, along with the transport method used to send this statement. The `transport` field can take two values: `silk` for the enhanced interconnect Postgres Pro Shardman mechanism, or `libpq` for sending via the standard PostgreSQL protocol.

14.6.4.1. Displaying Plans from the Remote Server

To see the execution plan that will be used on the remote server under the `EXPLAIN` block of the `ForeignScan` node, use the `postgres_fdw.foreign_explain` configuration parameter. The possible values are: `none` to exclude the `EXPLAIN` output from the remote servers, `full` to include the `EXPLAIN` output from the remote servers, `collapsed` to include the `EXPLAIN` output only for the first `ForeignScan` node under its `Append/MergeAppend`.

In production, it is recommended to disable this parameter (set it to `none`) or set it to `collapsed`, because obtaining any `EXPLAIN` information results in an additional implicit request to the server. Moreover, this request is executed in a synchronous mode, meaning the overall `EXPLAIN` output is built only once all the servers are sequentially queried. It can be a costly operation in case of a table with a large number of partitions.

Note that in case of the internal request for obtaining the `EXPLAIN` blocks for a remote plan, certain parameters are forcibly disabled, regardless of the parameters specified by a user when requesting `EXPLAIN` from the coordinator: `ANALYZE OFF`, `TIMING OFF`, `SUMMARY OFF`, `SETTINGS OFF`, `NETWORK OFF`. In this case, the `EXPLAIN` block of a remote plan will lack the corresponding metrics. Other `EXPLAIN` parameters (`FORMAT`, `VERBOSE`, `COSTS`, `BUFFERS`, `WAL`) are inherited from the coordinator.

The internal `EXPLAIN` request to a remote server is made with the `GENERIC_PLAN` option. Therefore, when analyzing the plans from the remote side, note that the `EXPLAIN` blocks display the generic plan.

14.6.4.2. Network Metrics and Latency

Setting the `NETWORK` option of the `EXPLAIN` command to `on` shows the network operation metrics for the plan nodes, including individual `ForeignScan` nodes and general nodes `Append` or `MergeAppend`.

For each plan node, the `FDW bytes`, `sent`, and `received` parameters are displayed for the outgoing and incoming traffic when the node is executed (regardless of the transport type). Note that these metrics are only output when the `ANALYZE` option of the `EXPLAIN` command is set to `on`.

When the `track_fdw_wait_timing` configuration parameter is enabled, the `wait_time` metric is also output. This metric summarizes all stages of the plan node execution, starting from the time the request is sent to the remote server, including the time spent on the execution itself and all the time until the complete set of results for that plan node is received.

Note that the `ForeignScan` node can operate in both synchronous and asynchronous modes. For the asynchronous execution, the node's execution function sends a request to the remote server and completes its execution without waiting for the result. The result is considered and processed later, upon receipt. In this scenario, the `wait_time` metric may not accurately reflect the actual execution time.

14.6.4.3. Query Tracing for Silk Transport

For the Silk transport, there is an option to output the extended debug information about tracing of a query passing from the coordinator to the remote server and back, including the results from the remote server. This information is only available if the `ANALYZE` option of the `EXPLAIN` command is set to `on`, and the `shardman.silk_tracepoints` configuration parameter is enabled.

When these parameters are enabled, each message transferred through the Silk transport (sending the SQL query, delivering it to the recipient, executing the query, and returning the execution result) is accompanied by an array of the timestamps measured at certain points in the pipeline. Once the query is executed, this information is displayed in the `EXPLAIN` block as rows starting with the word `Trace`. Each metric represents the difference between the timestamps at different points, in milliseconds:

Table 14.1. Query Tracing for Silk Transport Metrics

Interval	Description
bk shm->mp1 (qry)	The time taken to transfer an SQL query from the coordinator to its multiplexer via the shared memory.
mp1 shm->net (qry)	The time between receiving a query within the multiplexer from the shared memory and transferring it over the network.
net (qry)	The time spent by an SQL query to transfer over the network between the multiplexers.
mp2 recv->shm (qry)	The time between receiving an SQL query from the network and placing it in the queue in the shared memory on a remote multiplexer.
wk exec (1st tup)	The time spent to execute a query in <code>Silkworm</code> until the first row of the result is received.
wk exec (all tups)	The time spent to execute a query on <code>Silkworm</code> until the complete result is received.
wk->shm (1st tup)	The time taken to place the first row of the result into the <code>Silkworm</code> queue.
wk->shm (last tup)	The time taken to place the last row of the result into the <code>Silkworm</code> queue.
mp2 shm->net (1st tup)	The time between reading the first row of the result from the queue by the remote multiplexer and transferring it over the network.
net (1st tup)	The time spent to transfer the first row of the result over the network between the multiplexers.
mp1 recv->shm (1st tup)	The time between receiving the first row of the result from the network and placing it in the queue by the local multiplexer.
mp1 shm->bk (1st tup)	The time spent to retrieve the first row of the result from the queue by the coordinator.
mp2 shm->net (last tup)	The time between reading of the last row of the result from the queue by the remote multiplexer and transferring it over the network.
net (last tup)	The time spent to transfer the last row of the result over the network between the multiplexers.
mp1 recv->shm (last tup)	The time between receiving the last row of the result from the network and placing it in the queue by the local multiplexer.

Interval	Description
mp1 shm->bk (last tup)	The time taken by the coordinator to retrieve the last row o the result from the queue.
END-TO-END	The total time from sending the query to receiving the last row of the result. This approximately corresponds to the <code>wait_time</code> .

For the metrics `net (qry)`, `net (1st tup)`, and `net (last tup)`, the interval value is calculated as the difference between timestamps on different servers. Therefore, negative values may appear in these lines. This difference includes both time spent on a message transfer and the clock drift (positive or negative) between these servers. Thus, even with a slight drift, the values will be negative if its absolute value exceeds the duration of network transfer. Although it is not a bug, you should pay close attention to whether the cluster clocks are synchronized. For more information, see [Section 14.6.3](#).

Chapter 15. Parallel Query

Postgres Pro can devise query plans that can leverage multiple CPUs in order to answer queries faster. This feature is known as parallel query. Many queries cannot benefit from parallel query, either due to limitations of the current implementation or because there is no imaginable query plan that is any faster than the serial query plan. However, for queries that can benefit, the speedup from parallel query is often very significant. Many queries can run more than twice as fast when using parallel query, and some queries can run four times faster or even more. Queries that touch a large amount of data but return only a few rows to the user will typically benefit most. This chapter explains some details of how parallel query works and in which situations it can be used so that users who wish to make use of it can understand what to expect.

15.1. How Parallel Query Works

When the optimizer determines that parallel query is the fastest execution strategy for a particular query, it will create a query plan that includes a *Gather* or *Gather Merge* node. Here is a simple example:

```
EXPLAIN SELECT * FROM pgbench_accounts WHERE filler LIKE '%x%';
                                QUERY PLAN
-----
Gather  (cost=1000.00..217018.43 rows=1 width=97)
  Workers Planned: 2
    -> Parallel Seq Scan on pgbench_accounts  (cost=0.00..216018.33 rows=1 width=97)
        Filter: (filler ~~ '%x%':text)
(4 rows)
```

In all cases, the *Gather* or *Gather Merge* node will have exactly one child plan, which is the portion of the plan that will be executed in parallel. If the *Gather* or *Gather Merge* node is at the very top of the plan tree, then the entire query will execute in parallel. If it is somewhere else in the plan tree, then only the portion of the plan below it will run in parallel. In the example above, the query accesses only one table, so there is only one plan node other than the *Gather* node itself; since that plan node is a child of the *Gather* node, it will run in parallel.

Using [EXPLAIN](#), you can see the number of workers chosen by the planner. When the *Gather* node is reached during query execution, the process that is implementing the user's session will request a number of [background worker processes](#) equal to the number of workers chosen by the planner. The number of background workers that the planner will consider using is limited to at most [max_parallel_workers_per_gather](#). The total number of background workers that can exist at any one time is limited by both [max_worker_processes](#) and [max_parallel_workers](#). Therefore, it is possible for a parallel query to run with fewer workers than planned, or even with no workers at all. The optimal plan may depend on the number of workers that are available, so this can result in poor query performance. If this occurrence is frequent, consider increasing [max_worker_processes](#) and [max_parallel_workers](#) so that more workers can be run simultaneously or alternatively reducing [max_parallel_workers_per_gather](#) so that the planner requests fewer workers.

Every background worker process that is successfully started for a given parallel query will execute the parallel portion of the plan. The leader will also execute that portion of the plan, but it has an additional responsibility: it must also read all of the tuples generated by the workers. When the parallel portion of the plan generates only a small number of tuples, the leader will often behave very much like an additional worker, speeding up query execution. Conversely, when the parallel portion of the plan generates a large number of tuples, the leader may be almost entirely occupied with reading the tuples generated by the workers and performing any further processing steps that are required by plan nodes above the level of the *Gather* node or *Gather Merge* node. In such cases, the leader will do very little of the work of executing the parallel portion of the plan.

When the node at the top of the parallel portion of the plan is *Gather Merge* rather than *Gather*, it indicates that each process executing the parallel portion of the plan is producing tuples in sorted order, and that the leader is performing an order-preserving merge. In contrast, *Gather* reads tuples from the workers in whatever order is convenient, destroying any sort order that may have existed.

15.2. When Can Parallel Query Be Used?

There are several settings that can cause the query planner not to generate a parallel query plan under any circumstances. In order for any parallel query plans whatsoever to be generated, the following settings must be configured as indicated.

- `max_parallel_workers_per_gather` must be set to a value that is greater than zero. This is a special case of the more general principle that no more workers should be used than the number configured via `max_parallel_workers_per_gather`.

In addition, the system must not be running in single-user mode. Since the entire database system is running as a single process in this situation, no background workers will be available.

Even when it is in general possible for parallel query plans to be generated, the planner will not generate them for a given query if any of the following are true:

- The query writes any data or locks any database rows. If a query contains a data-modifying operation either at the top level or within a CTE, no parallel plans for that query will be generated. As an exception, the following commands, which create a new table and populate it, can use a parallel plan for the underlying `SELECT` part of the query:
 - `CREATE TABLE ... AS`
 - `SELECT INTO`
 - `CREATE MATERIALIZED VIEW`
 - `REFRESH MATERIALIZED VIEW`
- The query might be suspended during execution. In any situation in which the system thinks that partial or incremental execution might occur, no parallel plan is generated. For example, a cursor created using `DECLARE CURSOR` will never use a parallel plan. Similarly, a PL/pgSQL loop of the form `FOR x IN query LOOP .. END LOOP` will never use a parallel plan, because the parallel query system is unable to verify that the code in the loop is safe to execute while parallel query is active.
- The query uses any function marked `PARALLEL UNSAFE`. Most system-defined functions are `PARALLEL SAFE`, but user-defined functions are marked `PARALLEL UNSAFE` by default. See the discussion of [Section 15.4](#).
- The query is running inside of another query that is already parallel. For example, if a function called by a parallel query issues an SQL query itself, that query will never use a parallel plan. This is a limitation of the current implementation, but it may not be desirable to remove this limitation, since it could result in a single query using a very large number of processes.

Even when a parallel query plan is generated for a particular query, there are several circumstances under which it will be impossible to execute that plan in parallel at execution time. If this occurs, the leader will execute the portion of the plan below the `Gather` node entirely by itself, almost as if the `Gather` node were not present. This will happen if any of the following conditions are met:

- No background workers can be obtained because of the limitation that the total number of background workers cannot exceed `max_worker_processes`.
- No background workers can be obtained because of the limitation that the total number of background workers launched for purposes of parallel query cannot exceed `max_parallel_workers`.
- The client sends an `Execute` message with a non-zero fetch count. See the discussion of the [extended query protocol](#). Since `libpq` currently provides no way to send such a message, this can only occur when using a client that does not rely on `libpq`. If this is a frequent occurrence, it may be a good idea to set `max_parallel_workers_per_gather` to zero in sessions where it is likely, so as to avoid generating query plans that may be suboptimal when run serially.

15.3. Parallel Plans

Because each worker executes the parallel portion of the plan to completion, it is not possible to simply take an ordinary query plan and run it using multiple workers. Each worker would produce a full copy

of the output result set, so the query would not run any faster than normal but would produce incorrect results. Instead, the parallel portion of the plan must be what is known internally to the query optimizer as a *partial plan*; that is, it must be constructed so that each process that executes the plan will generate only a subset of the output rows in such a way that each required output row is guaranteed to be generated by exactly one of the cooperating processes. Generally, this means that the scan on the driving table of the query must be a parallel-aware scan.

15.3.1. Parallel Scans

The following types of parallel-aware table scans are currently supported.

- In a *parallel sequential scan*, the table's blocks will be divided into ranges and shared among the cooperating processes. Each worker process will complete the scanning of its given range of blocks before requesting an additional range of blocks.
- In a *parallel bitmap heap scan*, one process is chosen as the leader. That process performs a scan of one or more indexes and builds a bitmap indicating which table blocks need to be visited. These blocks are then divided among the cooperating processes as in a parallel sequential scan. In other words, the heap scan is performed in parallel, but the underlying index scan is not.
- In a *parallel index scan* or *parallel index-only scan*, the cooperating processes take turns reading data from the index. Currently, parallel index scans are supported only for btree indexes. Each process will claim a single index block and will scan and return all tuples referenced by that block; other processes can at the same time be returning tuples from a different index block. The results of a parallel btree scan are returned in sorted order within each worker process.

Other scan types, such as scans of non-btree indexes, may support parallel scans in the future.

15.3.2. Parallel Joins

Just as in a non-parallel plan, the driving table may be joined to one or more other tables using a nested loop, hash join, or merge join. The inner side of the join may be any kind of non-parallel plan that is otherwise supported by the planner provided that it is safe to run within a parallel worker. Depending on the join type, the inner side may also be a parallel plan.

- In a *nested loop join*, the inner side is always non-parallel. Although it is executed in full, this is efficient if the inner side is an index scan, because the outer tuples and thus the loops that look up values in the index are divided over the cooperating processes.
- In a *merge join*, the inner side is always a non-parallel plan and therefore executed in full. This may be inefficient, especially if a sort must be performed, because the work and resulting data are duplicated in every cooperating process.
- In a *hash join* (without the "parallel" prefix), the inner side is executed in full by every cooperating process to build identical copies of the hash table. This may be inefficient if the hash table is large or the plan is expensive. In a *parallel hash join*, the inner side is a *parallel hash* that divides the work of building a shared hash table over the cooperating processes.

15.3.3. Parallel Aggregation

Postgres Pro supports parallel aggregation by aggregating in two stages. First, each process participating in the parallel portion of the query performs an aggregation step, producing a partial result for each group of which that process is aware. This is reflected in the plan as a `Partial Aggregate` node. Second, the partial results are transferred to the leader via `Gather` or `Gather Merge`. Finally, the leader re-aggregates the results across all workers in order to produce the final result. This is reflected in the plan as a `Finalize Aggregate` node.

Because the `Finalize Aggregate` node runs on the leader process, queries that produce a relatively large number of groups in comparison to the number of input rows will appear less favorable to the query planner. For example, in the worst-case scenario the number of groups seen by the `Finalize Aggregate` node could be as many as the number of input rows that were seen by all worker processes in the `Partial Aggregate` stage. For such cases, there is clearly going to be no performance benefit to

using parallel aggregation. The query planner takes this into account during the planning process and is unlikely to choose parallel aggregate in this scenario.

Parallel aggregation is not supported in all situations. Each aggregate must be [safe](#) for parallelism and must have a combine function. If the aggregate has a transition state of type `internal`, it must have serialization and deserialization functions. See [CREATE AGGREGATE](#) for more details. Parallel aggregation is not supported if any aggregate function call contains `DISTINCT` or `ORDER BY` clause and is also not supported for ordered set aggregates or when the query involves `GROUPING SETS`. It can only be used when all joins involved in the query are also part of the parallel portion of the plan.

15.3.4. Parallel Append

Whenever Postgres Pro needs to combine rows from multiple sources into a single result set, it uses an `Append` or `MergeAppend` plan node. This commonly happens when implementing `UNION ALL` or when scanning a partitioned table. Such nodes can be used in parallel plans just as they can in any other plan. However, in a parallel plan, the planner may instead use a `Parallel Append` node.

When an `Append` node is used in a parallel plan, each process will execute the child plans in the order in which they appear, so that all participating processes cooperate to execute the first child plan until it is complete and then move to the second plan at around the same time. When a `Parallel Append` is used instead, the executor will instead spread out the participating processes as evenly as possible across its child plans, so that multiple child plans are executed simultaneously. This avoids contention, and also avoids paying the startup cost of a child plan in those processes that never execute it.

Also, unlike a regular `Append` node, which can only have partial children when used within a parallel plan, a `Parallel Append` node can have both partial and non-partial child plans. Non-partial children will be scanned by only a single process, since scanning them more than once would produce duplicate results. Plans that involve appending multiple results sets can therefore achieve coarse-grained parallelism even when efficient partial plans are not available. For example, consider a query against a partitioned table that can only be implemented efficiently by using an index that does not support parallel scans. The planner might choose a `Parallel Append` of regular `Index Scan` plans; each individual index scan would have to be executed to completion by a single process, but different scans could be performed at the same time by different processes.

[enable_parallel_append](#) can be used to disable this feature.

15.3.5. Parallel Plan Tips

If a query that is expected to do so does not produce a parallel plan, you can try reducing [parallel_setup_cost](#) or [parallel_tuple_cost](#). Of course, this plan may turn out to be slower than the serial plan that the planner preferred, but this will not always be the case. If you don't get a parallel plan even with very small values of these settings (e.g., after setting them both to zero), there may be some reason why the query planner is unable to generate a parallel plan for your query. See [Section 15.2](#) and [Section 15.4](#) for information on why this may be the case.

When executing a parallel plan, you can use `EXPLAIN (ANALYZE, VERBOSE)` to display per-worker statistics for each plan node. This may be useful in determining whether the work is being evenly distributed between all plan nodes and more generally in understanding the performance characteristics of the plan.

15.4. Parallel Safety

The planner classifies operations involved in a query as either *parallel safe*, *parallel restricted*, or *parallel unsafe*. A parallel safe operation is one that does not conflict with the use of parallel query. A parallel restricted operation is one that cannot be performed in a parallel worker, but that can be performed in the leader while parallel query is in use. Therefore, parallel restricted operations can never occur below a `Gather` or `Gather Merge` node, but can occur elsewhere in a plan that contains such a node. A parallel unsafe operation is one that cannot be performed while parallel query is in use, not even in the leader. When a query contains anything that is parallel unsafe, parallel query is completely disabled for that query.

The following operations are always parallel restricted:

- Scans of common table expressions (CTEs).
- Scans of temporary tables.
- Scans of foreign tables, unless the foreign data wrapper has an `IsForeignScanParallelSafe` API that indicates otherwise.
- Plan nodes that reference a correlated `SubPlan`.

15.4.1. Parallel Labeling for Functions and Aggregates

The planner cannot automatically determine whether a user-defined function or aggregate is parallel safe, parallel restricted, or parallel unsafe, because this would require predicting every operation that the function could possibly perform. In general, this is equivalent to the Halting Problem and therefore impossible. Even for simple functions where it could conceivably be done, we do not try, since this would be expensive and error-prone. Instead, all user-defined functions are assumed to be parallel unsafe unless otherwise marked. When using `CREATE FUNCTION` or `ALTER FUNCTION`, markings can be set by specifying `PARALLEL SAFE`, `PARALLEL RESTRICTED`, or `PARALLEL UNSAFE` as appropriate. When using `CREATE AGGREGATE`, the `PARALLEL` option can be specified with `SAFE`, `RESTRICTED`, or `UNSAFE` as the corresponding value.

Functions and aggregates must be marked `PARALLEL UNSAFE` if they write to the database, change the transaction state (other than by using a subtransaction for error recovery), access sequences, or make persistent changes to settings. Similarly, functions must be marked `PARALLEL RESTRICTED` if they access temporary tables, client connection state, cursors, prepared statements, or miscellaneous backend-local state that the system cannot synchronize across workers. For example, `setseed` and `random` are parallel restricted for this last reason.

In general, if a function is labeled as being safe when it is restricted or unsafe, or if it is labeled as being restricted when it is in fact unsafe, it may throw errors or produce wrong answers when used in a parallel query. C-language functions could in theory exhibit totally undefined behavior if mislabeled, since there is no way for the system to protect itself against arbitrary C code, but in most likely cases the result will be no worse than for any other function. If in doubt, it is probably best to label functions as `UNSAFE`.

If a function executed within a parallel worker acquires locks that are not held by the leader, for example by querying a table not referenced in the query, those locks will be released at worker exit, not end of transaction. If you write a function that does this, and this behavior difference is important to you, mark such functions as `PARALLEL RESTRICTED` to ensure that they execute only in the leader.

Note that the query planner does not consider deferring the evaluation of parallel-restricted functions or aggregates involved in the query in order to obtain a superior plan. So, for example, if a `WHERE` clause applied to a particular table is parallel restricted, the query planner will not consider performing a scan of that table in the parallel portion of a plan. In some cases, it would be possible (and perhaps even efficient) to include the scan of that table in the parallel portion of the query and defer the evaluation of the `WHERE` clause so that it happens above the `Gather` node. However, the planner does not do this.

Chapter 16. Autonomous Transactions

An *autonomous transaction* is an independent transaction run within a parent transaction. Unlike subtransactions, which can only be committed as part of the transaction they belong to, autonomous transactions must be committed or rolled back before their parent is finished. While subtransactions are mostly used for error handling and in stored procedures, the main purpose of autonomous transactions is to implement audits, when an attempt to perform a transaction must be logged regardless of whether this transaction has been committed successfully.

Note

Behavior of an autonomous transaction differs from the behavior of an ordinary transaction or a subtransaction. It is recommended to read the [Section 16.1](#) and [Section 16.2](#) sections before working with autonomous transactions.

16.1. Behavior

Although an autonomous transaction is always run inside another transaction, only a single transaction can be active at any given moment; other transactions have to wait until the active transaction is complete. When an autonomous transaction is started, its parent transaction is paused and pushed into the autonomous transaction stack. Once the autonomous transaction is committed or rolled back, the parent transaction is popped from the stack and resumed. You cannot commit the parent transaction until all its autonomous transactions are complete.

Within a single parent transaction, you can start several autonomous ones, which can be both consecutive and nested. By default, Postgres Pro allows running up to 100 autonomous transactions in all sessions simultaneously. You can increase this limit by changing the [max_autonomous_transactions](#) configuration parameter. The maximum nesting level is restricted to 128 and cannot be changed.

Autonomous transactions do not support the following scenarios, so they will cause an error if attempted:

- The `FOR UPDATE` locking clause is not supported, so you cannot lock a row in a parent transaction and skip it in the autonomous transaction.
- Temporary namespaces cannot be accessed from autonomous transactions.
- You cannot redefine isolation level for autonomous transactions within PL/pgSQL or PL/Python blocks.
- Logical replication slots cannot be used in autonomous transactions.
- Using foreign data wrappers that rely on PostgreSQL transaction mechanism in autonomous transactions is not supported, except for the adapted `postgres_fdw`. Doing so is likely to cause errors, such as issues with error handling during transaction rollback.
- Autonomous transactions cannot be used together with the following extensions:
 - `in_memory`
 - `online_analyze`

Consider the following examples. A continuous line denotes an active transaction, while a dotted line denotes a transaction which has been paused and pushed into the autonomous transaction stack. Time flows downwards.

```
BEGIN; -- starts ordinary transaction T0
|
INSERT INTO t VALUES (1);
:\
: BEGIN AUTONOMOUS TRANSACTION; -- starts autonomous transaction
```

```

: |                                -- T1, pushes T0 into stack
: |
: INSERT INTO t VALUES (2);
: |
: COMMIT AUTONOMOUS TRANSACTION / ROLLBACK AUTONOMOUS TRANSACTION;
: |                                -- ends autonomous transaction
: |                                -- T1, pops transaction T0 from stack
:/
COMMIT / ROLLBACK;                -- ends transaction T0

```

Depending on the two choices between `COMMIT` and `ROLLBACK`, the following query can return four different results:

```
SELECT sum(x) FROM t;
```

A parent transaction can have several autonomous transactions, which can be nested or follow one another:

```

BEGIN;                            -- starts ordinary transaction T0
|
INSERT INTO t VALUES (1);
:\
: BEGIN AUTONOMOUS TRANSACTION;    -- starts autonomous transaction
: |                                -- T1, pushes T0 into stack
: |
: INSERT INTO t VALUES (2);
: |
: COMMIT AUTONOMOUS TRANSACTION / ROLLBACK AUTONOMOUS TRANSACTION;
: |                                -- ends autonomous transaction
: |                                -- T1, pops T0 from stack
:/
|
:\
: BEGIN AUTONOMOUS TRANSACTION;    -- starts autonomous transaction
: |                                -- T2, pushes T0 into stack
: |
: INSERT INTO t VALUES (4);
: :\
: : BEGIN AUTONOMOUS TRANSACTION;  -- starts autonomous transaction
: : |                                -- T3, pushes T2 into stack
: : |
: : INSERT INTO t VALUES (6);
: : |
: : COMMIT AUTONOMOUS TRANSACTION / ROLLBACK AUTONOMOUS TRANSACTION;
: : |                                -- ends autonomous transaction
: : |                                -- T3, pops T2 from stack
: : /
: |
: COMMIT AUTONOMOUS TRANSACTION / ROLLBACK AUTONOMOUS TRANSACTION;
: |                                -- ends autonomous transaction
: |                                -- T2, pops T0 from stack
:/
COMMIT / ROLLBACK;                -- ends transaction T0

```

16.2. Visibility

Visibility rules for autonomous transactions are the same as for independent transactions executed via `dblink`. Autonomous transactions do not see the effects of their parent transaction because the latter has not been committed yet. At the same time, if other autonomous transactions precede the active au-

onomous transaction and are already committed, then the active transaction will see the effects of these transactions. The parent transaction might see the effects of its autonomous transactions, depending on its own [isolation level](#). For example, if the isolation level is Read Committed, the parent transaction will see all changes made by autonomous transactions. However, if the parent transaction is using the Repeatable Read isolation level, the changes made by autonomous transactions will be invisible.

Consider the following example, showing a set of both consequent and nested transactions. Each transaction in the set is labeled with its status. In addition, the transactions that have been already committed are highlighted.

```
- Ordinary transaction T0                [PAUSED]
  - Autonomous transaction A0            [COMMITTED]
  - Autonomous transaction A1            [COMMITTED]
  - Autonomous transaction A2            [PAUSED]
    - Autonomous transaction A2.0        [COMMITTED]
    - Autonomous transaction A2.1        [PAUSED]
      - Autonomous transaction A2.1.0    [ACTIVE]
```

In the above example, the active autonomous transaction A2.1.0 does not see the effects of the preceding transactions T0, A2, and A2.1 that are paused and not committed yet, regardless of their type. However, the active transaction sees the effects of the preceding autonomous transactions A0, A1, and A2.0 that have been already committed.

Autonomous transactions can cause single-session deadlocks, since an autonomous transaction can conflict with one of the paused transactions in its session. If an autonomous transaction tries to obtain any resource locked by its parent, a deadlock is reported.

16.3. SQL Grammar Extension for Autonomous Transactions

In Postgres Pro Shardman, [BEGIN/START TRANSACTION](#), [COMMIT/END](#), and [ROLLBACK/ABORT](#) transaction statements are extended with an optional keyword `AUTONOMOUS`. For example:

```
BEGIN [AUTONOMOUS] [TRANSACTION] [ ISOLATION LEVEL isolation-level]
END [AUTONOMOUS] [TRANSACTION]
```

where *isolation-level* can be one of the following: `SERIALIZABLE`, `REPEATABLE READ`, `READ COMMITTED`, or `READ UNCOMMITTED`.

The `AUTONOMOUS` keyword is required to start an autonomous transaction, but can be omitted when completing an autonomous transaction. If you have several nested levels of autonomous transactions, they are completed one by one, starting from the innermost nesting level. The top-level transaction cannot be autonomous.

16.4. PL/pgSQL Grammar Extension for Autonomous Transactions

Block construction in PL/pgSQL is extended by the optional `autonomous` keyword. It is possible to treat the whole function body as an autonomous transaction:

```
CREATE FUNCTION foo(x integer) RETURNS integer AS $$
BEGIN AUTONOMOUS
  RETURN x;
END;
$$ LANGUAGE plpgsql;
```

or create a separate `BEGIN/END` block:

```
CREATE OR REPLACE FUNCTION myaudit() RETURNS boolean AS $$
BEGIN AUTONOMOUS
```

```
BEGIN AUTONOMOUS
INSERT INTO audit_schedule VALUES ('new audit',now());
END;
... -- do audit itself
RETURN true;
END;
$$ LANGUAGE plpgsql;
```

Note

You cannot redefine isolation level for autonomous transactions within PL/pgSQL blocks.

If an exception is raised inside a `BEGIN AUTONOMOUS` block, this autonomous transaction is aborted, and the standard exception handling procedure is started, unwinding the stack and executing exception handlers until the exception is caught. Thus, exception handling for autonomous transactions in PL/pgSQL is done in the same way as for Postgres Pro subtransactions.

If an error is caught by an `EXCEPTION` clause, local variables of the PL/pgSQL function remain as they were when the error occurred, but all changes to persistent database state within the block are rolled back.

16.5. PL/Python Extension for Autonomous Transactions

In addition to the `subtransaction` method, PL/Python module provides the `autonomous` method that can be used in the `WITH` clause to start an autonomous transaction:

```
CREATE OR REPLACE FUNCTION pythonomous() RETURNS void AS $$
    plpy.execute("INSERT INTO atx_test VALUES ('asd', 123)")

    try:
        with plpy.autonomous():
            plpy.execute("INSERT INTO atx_test VALUES ('bsd', 456)")
    except plpy.SPIError, e:
        print("error: %s" % e.args)

    plpy.execute("INSERT INTO atx_test VALUES ('csd', 'csd')")
$$ LANGUAGE plpythonu;
```

Note

You cannot redefine isolation level for autonomous transactions within PL/Python blocks.

Exception handling for autonomous transactions in PL/Python is done in the same way as for subtransactions.

Part III. Server Administration

This part covers topics that are of interest to a Postgres Pro administrator. This includes installation, configuration of the server, management of users and databases, and maintenance tasks. Anyone running Postgres Pro server, even for personal use, but especially in production, should be familiar with these topics.

The information attempts to be in the order in which a new user should read it. The chapters are self-contained and can be read individually as desired. The information is presented in a narrative form in topical units. Readers looking for a complete description of a command are encouraged to review the [Part VI](#).

The first few chapters are written so they can be understood without prerequisite knowledge, so new users who need to set up their own server can begin their exploration. The rest of this part is about tuning and management; that material assumes that the reader is familiar with the general use of the Postgres Pro database system. Readers are encouraged review the [Part I](#) and [Part II](#) parts for additional information.

Chapter 17. Binary Installation

Warning

Mind that this chapter describes how to install Postgres Pro Enterprise without Postgres Pro Shardman. Following these steps will irreversibly install Postgres Pro Enterprise only, thus making it impossible to install Postgres Pro Shardman. For Postgres Pro Shardman installation guide, refer to [Quickstart Guide](#).

17.1. Installing Postgres Pro Enterprise

Postgres Pro only runs on Linux-based operating systems, and it is shipped as binary packages. Each Postgres Pro binary distribution consists of several packages. The package structure differs from vanilla PostgreSQL and offers the following installation modes:

- **Quick installation and setup.** The `postgrespro-sdm-17` package installs and configures all the components required for a viable ready-to-use configuration of both server and client components. Choose this option if you are going to install a single Postgres Pro instance only, and you are not worried about possible conflicts with other PostgreSQL-based products.

Important

Installing the `postgrespro-sdm-17` package can delete existing installations of Postgres Pro and PostgreSQL-based products. Similarly, this Postgres Pro installation may be automatically removed if you later install another PostgreSQL-based product. Do not use this package for upgrades or migrations, or if you are going to maintain several installations on the same system.

- **Custom installation.** You can select any packages required for your purposes, including development packages. This option needs manual configuration, so a good grasp of Linux and understanding of PostgreSQL architecture are required. This is the only option to choose if you are going to use Postgres Pro in one of the following scenarios:
 - Install several Postgres Pro versions side by side, or together with other PostgreSQL-based products.
 - Perform an upgrade from a previous version, or migrate from a different PostgreSQL-based product.
 - Control Postgres Pro server execution using high availability software, such as pacemaker, instead of the standard system service management facility.

The minimum hardware required to install Postgres Pro, create a database cluster and start the database server are as follows:

- 1-GHz processor
- 1 GB of RAM
- 1 GB of disk space

Additional disk space is required for data or supporting components.

17.1.1. Supported Linux Distributions

Note

This is the last Postgres Pro release that supports Ubuntu 20.04. Starting from the next release, this operating system will no longer be supported.

Postgres Pro binary packages are available for the following Linux-based systems:

- Red Hat Enterprise Linux (RHEL) systems and its derivatives: CentOS 8, Rocky Linux 8/9, Red Hat Enterprise Linux 8/9, Oracle Linux 8/9, ROSA Chrome 2021.2, Red OS Murom 7.3/8, AlmaLinux 9
- Debian-based systems: Debian 11/12, Ubuntu 20.04/22.04/24.04, Astra Linux 1.7/1.8/3.8
- ALT 10/11, ALT SP Release 10
- SUSE Linux Enterprise Server (SLES) 15
- Linux systems for Elbrus CPU architecture: ALT 9/10 for e2kv3/e2kv4, ALT 8.2 SP for e2kv3/e2kv4, ALT SP Release 10, Astra Linux Leningrad 8.1
- Linux systems for ARM architecture: ALT 9/10/11, ALT SP 8.2 and higher, Astra Linux Novorossiysk, Debian 10/11, Red OS Murom 7.3/8, RHEL 8/9, Ubuntu 20.04/22.04/24.04

Note

Postgres Pro binary packages rely on the tzdata library provided by the operating system, so you must ensure that the latest available version is installed. If tzdata is outdated, the time in your database may be incorrect.

17.1.2. Adding Package Repository

Before installing Postgres Pro from binary packages using the package manager of your operating system, you need to add the Postgres Pro package repository.

To add the repository, complete the following steps:

1. Contact the Postgres Pro Shardman support team to download the `pgpro-repo-add.sh` Bash script. The support team will also provide you with the credentials to access the repository.

This script adds the repository so that the package manager of your operating system can use it.

2. Make sure that the `wget` utility is installed and can access the Internet.

The script requires `wget` and the Internet access to check the validity of the credentials when adding the repository.

Note

It is not possible to use a SOCKS5 proxy for accessing the Internet as `wget` does not support this type of proxy.

3. (Optional) Examine the content of the script, if necessary.
4. Run the script as `root` (e. g. using `sudo`).
5. When prompted, enter the access credentials.

Important

If no credentials are provided or they are invalid, the script will offer to add the repository without credentials.

If you proceed without credentials, you will be required to manually provide the credentials later in order to access the repository.

Now you are ready to install Postgres Pro Shardman using the package manager of your operating system.

17.1.3. Quick Installation and Setup

If you only need to install a single Postgres Pro instance and are not going to use any other PostgreSQL-based products on your system, you can use the quick installation mode. The typical process is as follows:

1. [Add the package repository required for your operating system.](#)
2. Install the `postgrespro-sdm-17` package. It will bring all the required components via dependencies, create the default database, start the database server, as well as enable server autostart at system boot and make all the provided programs available in `PATH`.

Note

By default, the database configuration is set for the Postgres Pro product being installed. If you need your database configured for a different product, choose [custom installation](#) and use the `tune` argument of the `pg-setup initdb` command.

Once the installation completes, you can launch `psql` on behalf of the `postgres` user and connect to the newly created database, which is located in the `/var/lib/pgpro/sdm-17/data` directory.

Since the default database is created using the `pg-setup` script, the path to its data directory is stored in the `/etc/default/postgrespro-sdm-17` file. All the subsequent `pg-setup` commands, as well as any commands that manage Postgres Pro service, affect this database only. In this file, you can also change the value of `PG_OOM_ADJUST_VALUE` for postmaster child processes (see [Section 18.4.4](#) for details).

17.1.4. Custom Installation

Splitting the distribution into multiple packages enables customizing the installation for different purposes: database servers, client systems, or development workstations. Custom installations need to be configured manually, but give you more flexibility in using the product. You can install several Postgres Pro versions side by side, as well as together with other PostgreSQL-based products. In particular, this may be required when performing upgrades, or migrating from a different PostgreSQL-based product.

To perform a custom installation, complete the following steps:

1. [Add the package repository required for your operating system.](#)
2. [Choose Postgres Pro packages](#) required for your purposes and install them using the standard installation commands for your Linux distribution. The available packages are listed in [Table 17.1](#).

As a result, all files get installed into the `/opt/pgpro/sdm-17` directory.

3. Run `pg-wrapper` as root to make the installed client and server programs available via `PATH` and add SQL man pages to the man page configuration file. This utility is provided in the `postgrespro-sdm-17-client` package.

```
/opt/pgpro/sdm-17/bin/pg-wrapper links update
```

For details on how to handle possible conflicts, see [pg-wrapper](#) description.

4. If you chose to install the `postgrespro-sdm-17-server` package, make sure to complete the following server setup:

- a. Create the default database by running the helper script `pg-setup` as root with the `initdb` option:

```
/opt/pgpro/sdm-17/bin/pg-setup initdb [--tune=conf] [initdb_options]
```

where:

- the `tune` command-line argument sets the database configuration.

- `initdb_options` are regular [initdb](#) options.

Note

By default, `pg-setup` initializes the database cluster with checksums enabled. If this is not what you expect, specify the `--no-data-checksums`.

The `pg-setup` script performs database administration operations as user `postgres`. If you do not specify any [initdb](#) options, the default database is created in the `/var/lib/pgpro/sdm-17/data` directory, using localization settings specified in the `LANG` environment variable for the current session. All the `LC_*` environment variables are ignored.

Since the default database is created using the `pg-setup` script, the path to its data directory is stored in the `/etc/default/postgrespro-sdm-17` file. All the subsequent `pg-setup` commands, as well as any commands that manage Postgres Pro service, affect this database only. In this file, you can also change the value of `PG_OOM_ADJUST_VALUE` for postmaster child processes (see [Section 18.4.4](#) for details).

- b. Start the server by running `pg-setup` as root, as follows:

```
/opt/pgpro/sdm-17/bin/pg-setup service start
```

Like vanilla PostgreSQL, Postgres Pro server runs on behalf of the `postgres` user.

Note

By default, automatic server startup is disabled, so you can manually control the database recovery after a system reboot. Optionally, you can configure the Postgres Pro server to start automatically. For details, see [Section 17.1.4.2](#).

17.1.4.1. Choosing the Packages to Install

The table below lists all the available Postgres Pro Shardman packages.

Table 17.1. Postgres Pro Shardman Packages

Package	Description
<code>postgrespro-sdm-17</code>	Top-level package that installs and configures Postgres Pro for server and client systems. Do not use this package for upgrades or migrations.
<code>postgrespro-sdm-17-client</code>	Standard client applications, such as <code>psql</code> or <code>pg_dump</code> .
<code>postgrespro-sdm-17-libs</code>	Shared libraries required to deploy client applications, including <code>libpq</code> ; runtime libraries for ECPG processor.

Important

Installing the `postgrespro-sdm-17` package can delete existing installations of Postgres Pro and PostgreSQL-based products. Similarly, this Postgres Pro installation may be automatically removed if you later install another PostgreSQL-based product.

Package	Description
postgrespro-sdm-17-server	Postgres Pro server and PL/pgSQL server-side programming language.
postgrespro-sdm-17-contrib	Additional extensions and programs deployable on database servers.
postgrespro-sdm-17-devel	Header files and libraries for developing client applications and server extensions. On Debian-based systems, this package is called <code>postgrespro-sdm-17-dev</code> .
postgrespro-sdm-17-plperl	Server-side programming language based on Perl (see Chapter 46).
postgrespro-sdm-17-plpython3	Server-side programming language based on Python 3 (see Chapter 47).
postgrespro-sdm-17-pltcl	Server-side programming language based on Tcl (see Chapter 45).
postgrespro-sdm-17-docs	Documentation (English).
postgrespro-sdm-17-docs-ru	Documentation (Russian).
postgrespro-sdm-17-test	Test scripts for the server. This package is only available on RHEL-based and SUSE systems.
postgrespro-sdm-17-jit	This package provides support for Just-in-Time (JIT) compilation. This package is available only for x86_64 architecture and only for the supported Debian and Ubuntu systems, Astra Linux 1.7/1.8, supported ALT systems, CentOS 8, Rocky Linux 8/9, Oracle Linux 8/9, AlmaLinux 9, SLES 15 and RHEL 8/9. To learn more about enabling and using JIT, see Chapter 30 .
mamonsu	
pg-portal-modify-sdm-17	An extension to modify Postgres Pro cursors.
pg-probackup-sdm-17	pg_probackup utility.
pgpro-controldata	pgpro_controldata application to display control information of a PostgreSQL/Postgres Pro database cluster and compatibility information for a cluster and/or server.
pgpro-stats-sdm-17	pgpro_stats extension that tracks execution statistics of SQL statements, calculates wait event statistics and provides other useful metrics.

Also, Postgres Pro includes libraries and utilities required for the Postgres Pro server and provided packages. These packages are only provided for the distributions that don't include the required versions of these libraries.

Table 17.2. Third-party libraries and utilities for Postgres Pro Shardman

Library/utility	Description
liblz4	A library for fast lossless data compression.

Library/utility	Description
freetds	A library that implements TDS protocol. It is used by the <code>tds-fdw</code> extension. It includes <code>freetds-devel</code> , <code>freetds-doc</code> and <code>freetds-libs</code> . This package is available only for RHEL 8/9.
libunwind	A stack unwinding library used by <code>crash_info</code> . This package is provided for RHEL 8.
libbacktrace	A backtrace library used by <code>crash_info</code> . This package is provided for Red OS 7.3/8 and RHEL 8/9.
libsybdb5	A library that implements TDS protocol. This package is provided for SLES 15.
perl-IO-Tty, perl-IPC-Run	Two Perl libraries used in the TAP test framework available in the <code>postgrespro-sdm-17-devel</code> package for extension developers and engineers. <code>perl-IO-Tty</code> is provided for RED OS Murom 7.3/8, RHEL 8/9, and <code>perl-IPC-Run</code> is provided for RED OS Murom 7.3/8, RHEL 8/9, SLES 15.

Besides, there are separate packages providing several external modules that have been pre-built for compatibility with Postgres Pro Shardman:

Table 17.3. Third-party Packages Built for Postgres Pro Shardman

Package	Description
odbc-postgrespro-sdm-17	An Open Database Connectivity (ODBC) driver for accessing database management systems (DBMS). <code>odbc-postgresql</code> can be used instead, yet in some cases it may conflict with the outdated libraries of the vanilla PostgreSQL that do not provide functionality available in newer libraries.
orafce-sdm-17	This package implements in Postgres Pro some of the functions from the Oracle database that are missing (or behaving differently).
pgbouncer	pgbouncer — a connection pooler for Postgres Pro.
pgvector-sdm-17	pgvector extension that provides vector similarity search for Postgres Pro.
pldebugger-sdm-17	A set of shared libraries that implement an API for debugging PL/pgSQL functions in Postgres Pro.
plv8-sdm-17	The shared library PLV8 that provides a Postgres Pro procedural language powered by V8 Javascript Engine. This package is only available for the supported Debian and Ubuntu systems, ALT 10/11, SLES 15, RHEL 8/9, and Red OS Murom 7.3.
zstd	Command-line utility for the <code>libzstd</code> library.

Additionally, Postgres Pro provides separate packages with debug information for some operating systems:

- On Debian-based systems, see the `postgrespro-sdm-17-dbg` package.

- On RHEL-based systems, see the `postgrespro-sdm-17-debuginfo` package.
- On ALT Linux systems, all packages containing binary files have the corresponding `-debuginfo` packages.

Server installations require at least the following packages:

- `postgrespro-sdm-17-server`
- `postgrespro-sdm-17-client`
- `postgrespro-sdm-17-libs`

To use additional Postgres Pro extensions, you must also install the `postgrespro-sdm-17-contrib` package. On Debian-based systems, `postgrespro-sdm-17-server` package depends on `postgrespro-sdm-17-contrib` package, so the latter must always be installed together with the server.

For client installations, it is usually enough to install the `postgrespro-sdm-17-client` and `postgrespro-sdm-17-libs` packages. If you use custom applications and do not need standard client utilities such as `psql`, you can install the `postgrespro-sdm-17-libs` package only.

Development workstations require at least the following packages:

- `postgrespro-sdm-17-libs`
- `postgrespro-sdm-17-devel/ postgrespro-sdm-17-dev`

You may also want to install and configure the server with a test database on development systems. For details on additional configuration that may be required, [Section 17.1.5](#).

17.1.4.2. Enabling Automatic Server Startup

If you are running a custom installation, automatic server startup is disabled by default. Once the default database is created, you can configure the server to start automatically upon system boot using service management solutions available in your operating system or third-party high-availability software. To facilitate this task, `postgrespro-sdm-17-server` package provides the [pg-setup](#) script, which is installed in the `/opt/pgpro/sdm-17/bin` directory.

To enable server autostart, run the [pg-setup](#) script with the following options:

```
pg-setup service enable
```

If required, you can disable server autostart using the same script:

```
pg-setup service disable
```

Alternatively, you can use system service management solutions directly by running the autostart scripts for SysV-style `init.d` and `systemd` provided in the `postgrespro-sdm-17-server` package. Depending on your Linux distribution, Postgres Pro supports different service management solutions:

Linux Distribution	Provided Scripts
RHEL-based systems, SLES 15	<code>systemd</code> unit file
Debian, Ubuntu, ALT 10/11	Both <code>systemd</code> unit file and SysV-style <code>init.d</code> script

To use `systemd` for automatic server startup, run the following command:

```
systemctl enable postgrespro-sdm-17
```

To use SysV-style `init.d` script:

- On ALT Linux systems, run the following command:

```
systemctl enable postgrespro-sdm-17
```

- On Debian systems, use `update-rc.d`. See the corresponding man page for details.

17.1.5. Setting up Development Workstations

While installing `postgrespro-sdm-17-libs` and `postgrespro-sdm-17-devel/postgrespro-sdm-17-dev` packages may be enough, it is usually convenient to have the server set up on the development system. For quick setup, you can install `postgrespro-sdm-17` package, which automatically configures the provided client and server programs and creates the default database. However, if you are going to use several PostgreSQL-based products simultaneously, follow the custom installation instructions in [Section 17.1.4](#).

To compile programs with Postgres Pro libraries using the `pg_config` utility shipped with Postgres Pro, make sure it appears before the path to other `pg_config` versions, if any. Note that on RHEL-based systems `pg_config` is not added to `PATH` automatically. If you do not have any other `pg_config` versions on your system, you can use [pg-wrapper](#) provided in the `postgrespro-sdm-17-client` package to create a symbolic link to `pg_config` in the standard binary directory.

To compile programs using `pkg-config` command, add the `/opt/pgpro/sdm-17/lib/pkgconfig/` path to the `PKG_CONFIG_PATH` environment variable.

If you would like to compile Postgres Pro extensions that support [JIT inlining](#), make sure to meet the following additional requirements:

- Install LLVM development package and Clang compiler. You must choose the packages of the same version that was used for the `postgrespro-sdm-17-jit` to be installed on the server. To determine the required version for the current Postgres Pro release, check the `CLANG` value in the `/opt/pgpro/sdm-17/lib/pgxs/src/Makefile.global` file.
- When running `make` or `make install` commands, specify the `with-llvm=yes` option to compile and install bitcode files for your extension. By default, bitcode compilation is disabled as it depends on Clang compiler availability.

17.1.5.1. Using Third-Party Programs with Postgres Pro

To use Postgres Pro server with a client program provided with a third-party product, you can install the version of PostgreSQL libraries that was used to compile this program. For example, if this program is provided with vanilla PostgreSQL, you may need to install the `libpq` or `postgresql-libs` packages available for your Linux distribution. In this case, the program may not be able to use some new features of Postgres Pro server, but it is probably not designed to use them anyway.

If you prefer to use Postgres Pro libraries with a third-party program, or would like to enable support for a new feature that does not require client application change, such as SCRAM authentication, you can recompile your program with Postgres Pro libraries.

Important

If the program is compiled with one version of `libpq` but used with another, its stable work cannot be guaranteed.

If you are creating `.rpm` or `.deb` packages for your program, it is recommended to do the following:

1. Add `/opt/pgpro/sdm-17/bin` to `PATH` inside your build scripts (`.spec` files or `debian/rules`).
2. Specify `postgrespro-sdm-17-dev` in the `BuildDepends` or `BuildRequires` tags for your program.

Thus, you can ensure that your package build process calls the right version of `pg_config` whenever the source package is rebuilt.

17.1.6. Configuring Multiple Postgres Pro Instances

To set up several Postgres Pro server instances with different data directories on Linux, do the following:

1. Install and configure Postgres Pro as explained in [Section 17.1.3](#) or [Section 17.1.4](#).
2. Once the first default database is created, run `initdb` specifying the path to a different data directory and any other parameters required to initialize another server instance.
3. Specify different ports for your server instances in the corresponding `postgresql.conf` files to avoid conflicts.
4. If required, configure automatic server startup, as follows:
 - a. Create a copy of `/etc/init.d/postgrespro-sdm-17` or `/lib/systemd/system/postgrespro-sdm-17.service` with a different name, specifying the path to the data directory.
 - b. Enable automatic server startup using the provided autostart scripts for your system service management facility instead of `pg-setup`, as described in [Section 17.1.4.2](#). Make sure to use the re-named copies of the scripts you created in step 1.

17.1.7. Antivirus Considerations

It is strongly recommended to avoid using antivirus software on systems where Postgres Pro is running because it may cause additional load on your environment and result in unexpected database behavior that would lead to performance and reliability issues. If you need to use antivirus software, make sure to exclude the following directories from virus scanning as they do not contain any executable files:

- PGDATA directory that stores main cluster data, usually located in `/var/lib/pgpro/sdm-17/data` unless you specified another directory in `initdb` options
- Paths to created [tablespaces](#)

17.2. Installing Additional Supplied Modules

Postgres Pro comes with a set of additional server extensions, or modules. On Linux, these extensions are provided in the `postgrespro-contrib` package. On Windows, these extensions are installed together with the server components.

Once you have the binary files installed, you have to enable additional extensions in the database in order to use them. In most cases, you only need to issue the `CREATE EXTENSION` command. However, some extensions also require shared libraries to be preloaded on server startup. If you want to use such extensions, you need to configure parameter

```
shared_preload_libraries = 'lib1, lib2, lib3'
```

in the `postgresql.conf` file of your Postgres Pro database instance and restart the server before executing the `CREATE EXTENSION` statement.

For the exact installation and configuration instructions for each particular extension, see the corresponding documentation under [Appendix F](#).

To get the list of extensions available in your Postgres Pro installation, you can view the `pg_available_extensions` system catalog.

17.2.1. Installing New Extensions in Certified Product Editions

Working in a highly secured environment brings some restrictions. With a superuser prohibited due to its unlimited access rights, you can opt for regular users (for example, a DBMS Administrator) to handle operations.

While trusted extensions can be installed and manipulated with non-superuser rights, other extensions require an elevated security level. Installation of some extensions in a highly secured environment, for example `pg_proaudit`, follows a special procedure.

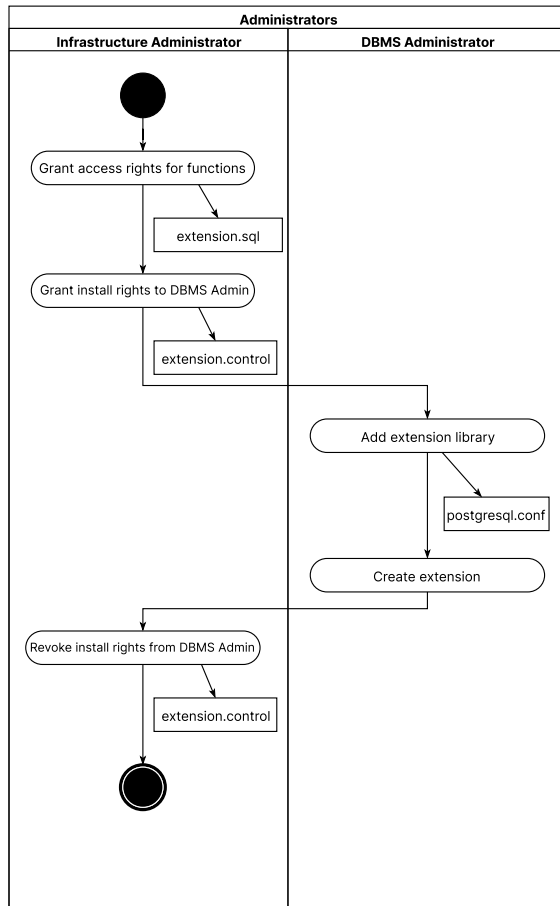
This procedure involves an Infrastructure Administrator and a DBMS Administrator:

- [An Infrastructure Administrator](#) is responsible for the overall system security and does not manage Postgres Pro.

- A **DBMS Administrator**, represented by the `PGPRO_DBMS_ADMIN` role in Postgres Pro, is responsible for configuration, setup, and administration of Postgres Pro DBMS.

The steps that the administrators need to take to install a new extension are shown in [Figure 17.1](#).

Figure 17.1. Steps taken by administrators to install an extension



17.2.1.1. Infrastructure Administrator: Steps to be Taken

Only the Infrastructure Administrator is allowed to modify configuration files of extensions in the `share/extension` directory.

The `.sql` file of an extension contains an SQL interface functions declaration. By default, these functions can be used only by a superuser. Under the extension installation procedure, the Infrastructure Administrator takes the following steps:

1. Adds the `GRANT` clauses to the `.sql` file, as shown below, to allow the DBMS Administrator to use these functions without requesting a superuser to grant such access:

```
-- Create new versions of objects
CREATE FUNCTION pg_proaudit_show()
RETURNS TABLE(db_name text,
               event_type text,
               object_type text,
               object_oid oid,
               role_name text)
AS 'MODULE_PATHNAME', 'pg_proaudit_show_conf'
LANGUAGE C VOLATILE;
REVOKE ALL ON FUNCTION pg_proaudit_show() FROM public;
```

The Infrastructure Administrator grants rights to a non-superuser:

```
-- Create new versions of objects
CREATE FUNCTION pg_proaudit_show()
RETURNS TABLE(db_name text,
               event_type text,
               object_type text,
               object_oid oid,
               role_name text)
AS 'MODULE_PATHNAME', 'pg_proaudit_show_conf'
LANGUAGE C VOLATILE;
REVOKE ALL ON FUNCTION pg_proaudit_show() FROM public;
GRANT ALL ON FUNCTION pg_proaudit_show() TO PGPRO_DBMS_ADMIN;
```

2. Allows a non-superuser to install the extension by changing or adding the `trusted` property in the `.control` file of the extension, thereby granting a temporary installation permission:

```
trusted = true
```

Enabling the usage of foreign data wrappers by the `PGPRO_DBMS_ADMIN` role requires a special security permission. To grant the permission, the Infrastructure Administrator adds the `GRANT USAGE ON FOREIGN DATA WRAPPER` command to a respective `.sql` file of the extension. Below is the example for [postgres_fdw](#):

```
GRANT USAGE ON FOREIGN DATA WRAPPER postgres_fdw TO PGPRO_DBMS_ADMIN;
```

17.2.1.2. DBMS Administrator: Steps to be Taken

The DBMS Administrator is allowed to modify Postgres Pro configuration files, except for `pg_hba.conf`, which stores security information. Only the Infrastructure Administrator is allowed to modify the `pg_hba.conf` configuration file. Under the extension installation procedure, the DBMS Administrator takes the following steps:

1. Adds a respective library file to the `shared_preload_libraries` variable of the `postgresql.conf` configuration file and reloads the database server for changes to take effect.
2. Creates the extension using the `CREATE EXTENSION` command.

17.2.1.3. Infrastructure Administrator: Final Step

To make further use of the extension secure, the Infrastructure Administrator reverts the `trusted` property to the original state:

- If it was not specified, it is deleted.
- If it was `FALSE`, it is changed back to `FALSE`.

17.3. Migration

17.3.1. Migrating to Postgres Pro Shardman

Different major versions of Postgres Pro Shardman, as well as different PostgreSQL-based products based on the same major version, can have binary incompatible databases, so you cannot replace the server binary and continue running.

For upgrade instructions specific to a particular release, see the [Release Notes](#) for the corresponding Postgres Pro Shardman version.

17.3.2. Migration of a Postgres Pro Shardman Database Schema

Let's use the demo database “Airlines” as an example for development. The detailed description of the database schema is available at <https://postgrespro.ru/education/demodb>. This schema is used as a demo in [training courses of Postgres Professional](#), for example, in “QPT. Query Optimization”.

The schema authors characterized it like this: “We tried to make the database schema as simple as possible, without overloading it with unnecessary details, but not too simple to allow writing interesting and meaningful queries.”

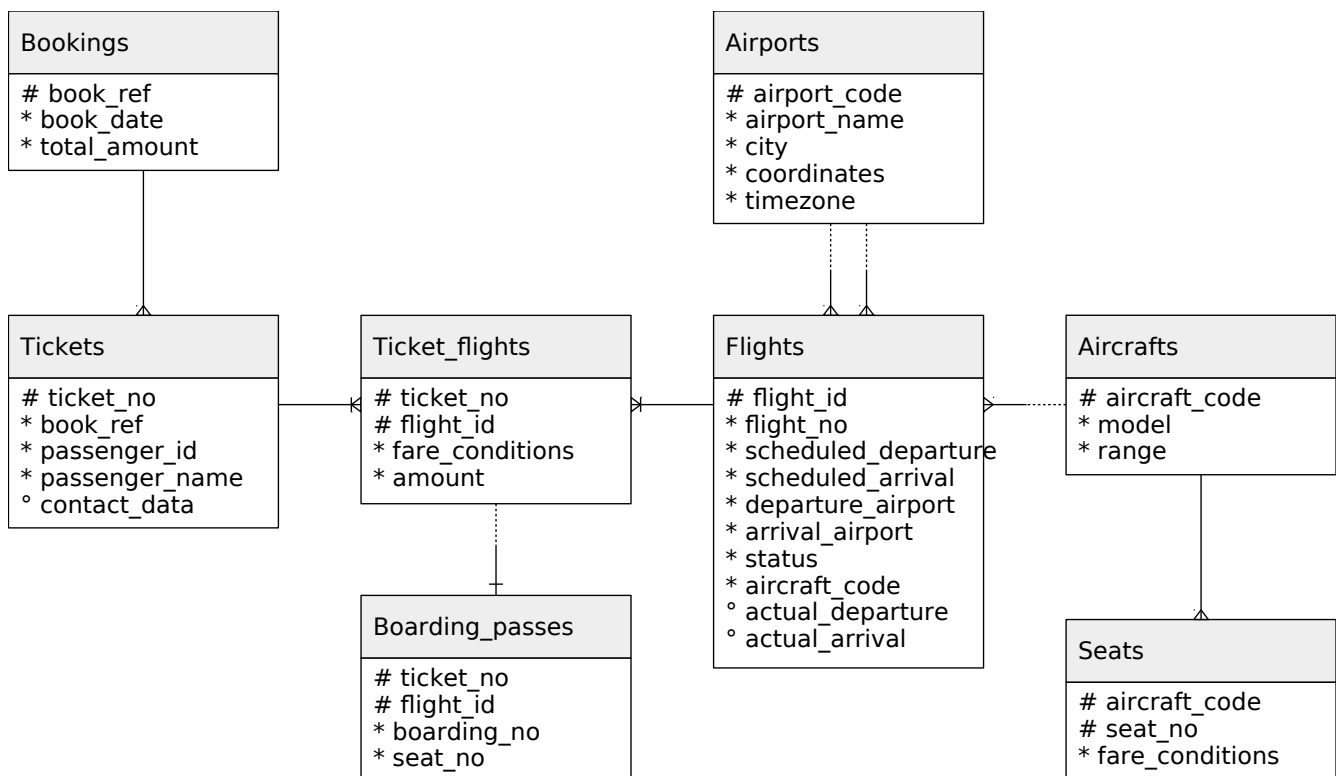
The database schema contains several tables with meaningful contents. For example, let's take the demo database version of 13.10.2016. You can find a link to downloading the database and schema dump (in Russian) following the link <https://postgrespro.ru/education/courses/QPT>. In addition to query examples provided below, you can find more examples from the above course and in the “*Postgres. The First Experience*” book.

This section shows two examples of schema modification and query adaptation:

- *Naive approach*. It is simple, with minimal transformations to the schema, and it aims to add clarity to how queries work in a distributed schema.
- *Complex approach*. It is more complex, provided for better understanding of problems and processes that a developer may confront when migrating to a distributed schema and adapting applications to such a schema.

17.3.2.1. Database Source Schema

Figure 17.2. Database Source Schema



The authors describe the “Airlines” database as follows:

The main entity is a booking (*bookings*).

One booking can include several passengers, with a separate ticket (*tickets*) issued to each passenger. A ticket has a unique number and includes information about the passenger. As such, the passenger is not a separate entity. Both the passenger's name and identity document number can change over time, so it is impossible to uniquely identify all the tickets of a particular person; for simplicity, we can assume that all passengers are unique.

The ticket includes one or more flight segments (*ticket_flights*). Several flight segments can be included into a single ticket if there are no non-stop flights between the

points of departure and destination (connecting flights), or if it is a round-trip ticket. Although there is no constraint in the schema, it is assumed that all tickets in the booking have the same flight segments.

Each flight (`flights`) goes from one airport (`airports`) to another. Flights with the same flight number have the same points of departure and destination, but differ in departure date.

At flight check-in, the passenger is issued a boarding pass (`boarding_passes`), where the seat number is specified. The passenger can check in for the flight only if this flight is included into the ticket. The flight-seat combination must be unique to avoid issuing two boarding passes for the same seat.

The number of seats (`seats`) in the aircraft and their distribution between different travel classes depends on the model of the aircraft (`aircrafts`) performing the flight. It is assumed that every aircraft model has only one cabin configuration. Database schema does not check that seat numbers in boarding passes have the corresponding seats in the aircraft (such verification can be done using table triggers, or at the application level).

Let's look at the common entities and sizes of tables in the above schema. It is clear that `ticket_flights`, `boarding_passes` and `tickets` tables are linked by the `ticket_no` field. Additionally, the data size in these tables is 95% the total DB size.

Let's look at the `bookings` table. Although it seems to have a pretty compact structure, it can reach a considerable size over time.

Migration examples are provided for a Postgres Pro Shardman cluster that contains four shards. Sharded tables are divided into four parts, so that one part of a sharded table is only located in one shard. This is done on purpose, to more clearly display query plans. In real life, the number of partitions should be determined by the maximum number of cluster nodes.

When migrating a real-life DB schema, you should think over in advance the number of partitions to partition data in distributed tables. Also bear in mind that the best migration approach is to use SQL transformations that impose minimal limitations on database objects.

17.3.2.2. Postgres Pro Shardman Cluster Configuration

The Postgres Pro Shardman cluster consists of four nodes — `node1`, `node2`, `node3` and `node4`. Each cluster node is a shard.

The examples assume that the tables are divided into four partitions by the sharding key (`num_parts = 4`) and distributed across cluster nodes. Each table part with the data is located in the corresponding shard:

- `shard-1` is located on the cluster node `node1`
- `shard-2` is located on the cluster node `node2`
- `shard-3` is located on the cluster node `node3`
- `shard-4` is located on the cluster node `node4`

The cluster is intentionally presented in a simplified configuration. Cluster nodes have no replicas, and the configuration is not fault-tolerant.

17.3.2.3. Selecting the Sharding Key

17.3.2.3.1. Naive¹ Approach — `ticket_no` Sharding Key

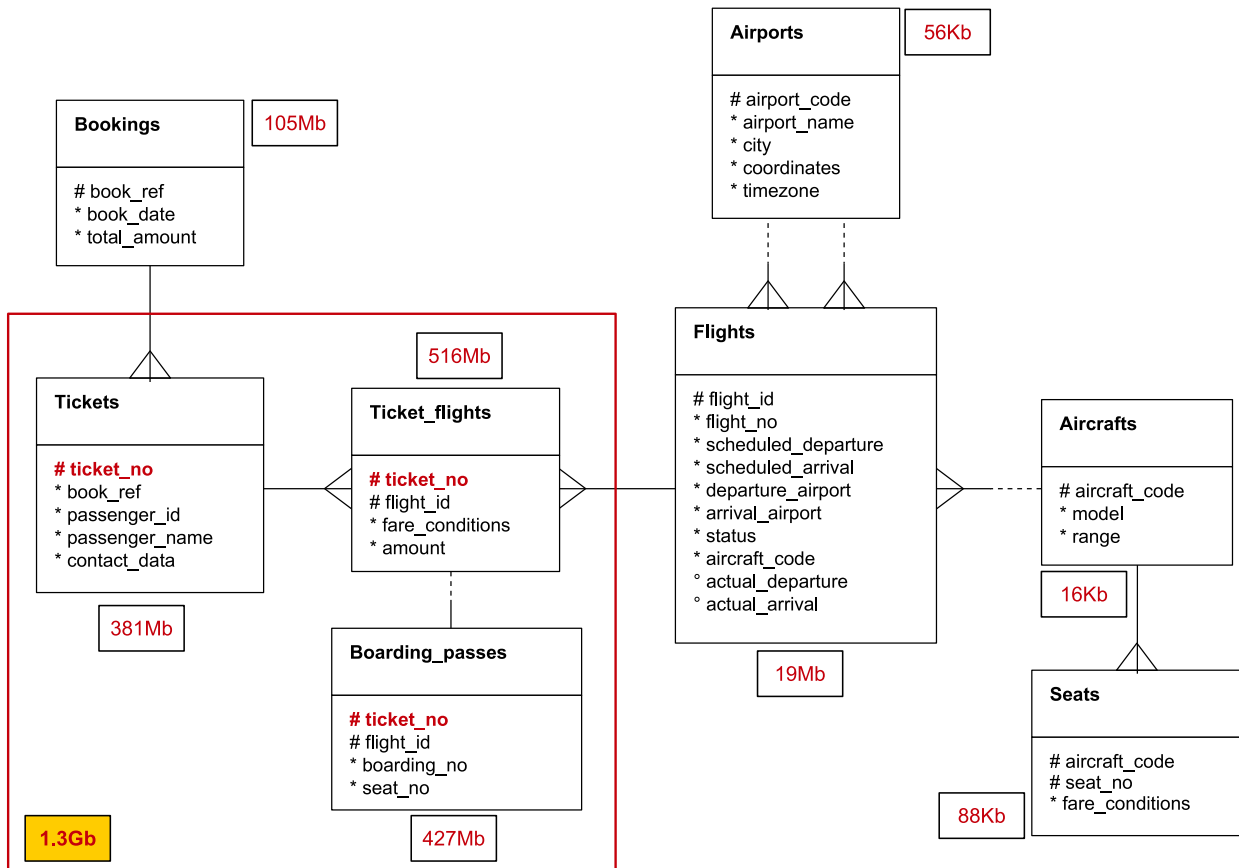
With this approach, the choice of the sharding key is pretty evident. It is the ticket number `ticket_no`. The ticket number is the primary key of the `tickets` table, and it is a foreign key of the `ticket_flights` and `boarding_passes` tables.

The primary key of the `ticket_flights` and `boarding_passes` tables is composite. It is a unique index composed of `ticket_no` and `flight_id`.

So if `ticket_no` is chosen to be a sharding key, the data of the three tables is distributed across cluster shards and partitions that contain linked data are located in the same shards.

The rest of the tables — `airports`, `flights`, `aircrafts` and `seats` are small enough and rarely change. This allows making them global tables, or dictionary tables.

Figure 17.3. Naive Approach Schema



The main advantage of this approach from the point of view of creating the schema and queries to the DB is that no changes are needed except those that are inherent to working with distributed systems, that is, explicitly declaring tables, sequences etc. as distributed when creating them.

Once the sharding key is selected, we can proceed to creation of the distributed schema.

1

17.3.2.3.1.1. Creating the Schema Distributed by `ticket_no`

First, turn on broadcasting DDL statements to all cluster shards:

```
SET shardman.broadcast_ddl TO on;
```

Let's create the `bookings` schema on all shards:

¹ In the context of computer science, the expression “naïve approach” (verbatim: naive method, naive approach) means something very similar to “brute-force approach” and means the first basic idea that occurs in one's mind and often takes no account of the complexity, corner cases and of some requirements. On one hand, this is a coarse and direct method that only aims to get a working solution. On the other hand, such solutions are easy to understand and implement, but system resources may be used inefficiently.

```
CREATE SCHEMA bookings;
```

As tables in the schema are linked with one another by a foreign key, the order of creating them, as well as auxiliary objects, matters.

The demo database contains “snapshots” of data, similar to a backup copy of a real system captured at some point in time. For example, if a flight has the `Departed` status, it means that the aircraft had already departed and was airborne at the time of the backup copy. The “snapshot” time is saved in the `bookings.now()` function. You can use this function in demo queries for cases where you would use the `now()` function in a real database. In addition, the return value of this function determines the version of the demo database. The latest version available is of 13.10.2016:

```
SELECT bookings.now() as now;
           now
-----
2016-10-13 17:00:00+03
```

In relation to this moment, all flights are classified as past and future flights.

Let's create the utility function `bookings.now()`:

```
CREATE FUNCTION bookings.now() RETURNS timestamp with time zone
LANGUAGE sql IMMUTABLE COST 0.009999999978
AS
$sql$
SELECT $qq$2016-10-13 17:00:00$qq$::TIMESTAMP AT TIME ZONE
$zz$Europe/Moscow$zz$;
$sql$;
```

In addition to tables, a global sequence is needed for generating IDs for data insertion in the `flights` table. In this example, we create the sequence explicitly and link it with a column of this table by assigning the generated values by default.

Let's create the sequence using the following DDL statement:

```
CREATE SEQUENCE bookings.flights_flight_id_seq
INCREMENT BY 1
NO MINVALUE
NO MAXVALUE
CACHE 1 with(global);
```

`with(global)` creates a single distributed sequence available on all cluster nodes, which assigns values in a certain range for each shard, and the ranges for different shards do not intersect. See [CREATE SEQUENCE](#) and [SQL Commands](#) for more details of global sequences.

Under the hood of global sequences, there are regular sequences on each shard, and they are allocated by sequential blocks (of 65536 numbers by default). When all the numbers in a block are over, the next block is allocated to the local sequence of the shard. I.e., numbers from the global sequences are unique, but there is no strict monotony, and there may be “holes” in the values given by the sequencer².

The sequences can have the `bigserial`, `smallserial`, or `serial` type. Sequences are applicable both for sharded and global tables.

You should not create local sequences in each shard as their values may be duplicated.

2

² As values from different ranges can be assigned, the value can leap. For example, the value of 5 may be assigned in the first shard, the value of 140003 — in the second one, 70003 — in the third one etc.

Now, we create global tables. As explained above, they are small-size, their data changes rarely, so they are actually dictionary tables, which must contain the same data in all cluster shards. It is required that each global table has a primary key.

Let's create global tables using the following DDL statements:

```
CREATE TABLE bookings.aircrafts (  
    aircraft_code character(3) NOT NULL primary key,  
    model text NOT NULL,  
    range integer NOT NULL,  
    CONSTRAINT aircrafts_range_check CHECK ((range > 0))  
) with (global);  
  
CREATE TABLE bookings.seats (  
    aircraft_code character(3) references bookings.aircrafts(aircraft_code),  
    seat_no character varying(4) NOT NULL,  
    fare_conditions character varying(10) NOT NULL,  
    CONSTRAINT seats_fare_conditions_check CHECK (((fare_conditions)::text = ANY  
(ARRAY[('Economy'::character varying)::text, ('Comfort'::character varying)::text,  
('Business'::character varying)::text]))),  
    PRIMARY KEY (aircraft_code, seat_no)  
) with (global);  
  
CREATE TABLE bookings.airports (  
    airport_code character(3) NOT NULL primary key,  
    airport_name text NOT NULL,  
    city text NOT NULL,  
    longitude double precision NOT NULL,  
    latitude double precision NOT NULL,  
    timezone text NOT NULL  
) with (global);  
  
CREATE TABLE bookings.bookings (  
    book_ref character(6) NOT NULL,  
    book_date timestamp with time zone NOT NULL,  
    total_amount numeric(10,2) NOT NULL,  
    PRIMARY KEY (book_ref)  
) with (global);  
  
CREATE TABLE bookings.flights (  
    flight_id bigint NOT NULL PRIMARY KEY, -- <= a sequence will be assigned  
    flight_no character(6) NOT NULL,  
    scheduled_departure timestamp with time zone NOT NULL,  
    scheduled_arrival timestamp with time zone NOT NULL,  
    departure_airport character(3) REFERENCES bookings.airports(airport_code),  
    arrival_airport character(3) REFERENCES bookings.airports(airport_code),  
    status character varying(20) NOT NULL,  
    aircraft_code character(3) references bookings.aircrafts(aircraft_code),  
    actual_departure timestamp with time zone,  
    actual_arrival timestamp with time zone,  
    CONSTRAINT flights_check CHECK ((scheduled_arrival > scheduled_departure)),  
    CONSTRAINT flights_check1 CHECK (((actual_arrival IS NULL) OR  
((actual_departure IS NOT NULL) AND (actual_arrival IS NOT NULL) AND (actual_arrival >  
actual_departure)))),  
    CONSTRAINT flights_status_check CHECK (((status)::text = ANY (ARRAY[('On  
Time'::character varying)::text, ('Delayed'::character varying)::text,  
('Departed'::character varying)::text, ('Arrived'::character varying)::text,  
('Scheduled'::character varying)::text, ('Cancelled'::character varying)::text])))  
) with (global);
```

```
-- associate the sequence with table column
ALTER SEQUENCE bookings.flights_flight_id_seq OWNED BY bookings.flights.flight_id;

-- assign the default value to the column
ALTER TABLE bookings.flights ALTER COLUMN flight_id SET DEFAULT
nextval('bookings.flights_flight_id_seq');

ALTER TABLE bookings.flights ADD CONSTRAINT
flights_flight_no_scheduled_departure_key UNIQUE (flight_no, scheduled_departure);
```

Next, we create sharded tables `tickets`, `ticket_flights` and `boarding_passes` in the `bookings` schema:

```
CREATE TABLE bookings.tickets (
    ticket_no character(13) PRIMARY KEY,
    book_ref character(6) REFERENCES bookings.bookings(book_ref),
    passenger_id character varying(20) NOT NULL,
    passenger_name text NOT NULL,
    contact_data jsonb
) with (distributed_by='ticket_no', num_parts=4);

CREATE TABLE bookings.ticket_flights (
    ticket_no character(13) NOT NULL,
    flight_id bigint references bookings.flights(flight_id),
    fare_conditions character varying(10) NOT NULL,
    amount numeric(10,2) NOT NULL,
    CONSTRAINT ticket_flights_amount_check CHECK ((amount >= (0)::numeric)),
    CONSTRAINT ticket_flights_fare_conditions_check CHECK (((fare_conditions)::text
= ANY (ARRAY[('Economy'::character varying)::text, ('Comfort'::character
varying)::text, ('Business'::character varying)::text]))),
    PRIMARY KEY (ticket_no, flight_id)
) with (distributed_by='ticket_no', colocate_with='bookings.tickets');

CREATE TABLE bookings.boarding_passes (
    ticket_no character(13) NOT NULL,
    flight_id bigint NOT NULL,
    boarding_no integer NOT NULL,
    seat_no character varying(4) NOT NULL,
    FOREIGN KEY (ticket_no, flight_id) REFERENCES
bookings.ticket_flights(ticket_no, flight_id),
    PRIMARY KEY (ticket_no, flight_id)
) with (distributed_by='ticket_no', colocate_with='bookings.tickets');

-- constraints must contain sharding key
ALTER TABLE bookings.boarding_passes ADD CONSTRAINT
boarding_passes_flight_id_boarding_no_key UNIQUE (ticket_no, flight_id, boarding_no);

ALTER TABLE bookings.boarding_passes ADD CONSTRAINT
boarding_passes_flight_id_seat_no_key UNIQUE (ticket_no, flight_id, seat_no);
```

Additionally, when creating sharded tables, the `num_parts` parameter can be specified, which defines the number of sharded table partitions. In this example, it equals 4 to minimize the output of query plans. The default value is 20. This parameter may be important if in future you are going to add shards to a cluster and scale horizontally.

Based on the assumed future load and data size, `num_parts` should be sufficient for data rebalancing when new shards are added (`num_parts` must be greater than or equal to the number of cluster nodes). On the other hand, too many partitions cause a considerable increase of the query planning time. Therefore, an optimal balance should be achieved between the number of partitions and number of cluster nodes.

The last thing to do is to create a view that is needed to execute some queries:

```
CREATE VIEW bookings.flights_v AS
SELECT f.flight_id,
       f.flight_no,
       f.scheduled_departure,
       timezone(dep.timezone, f.scheduled_departure) AS scheduled_departure_local,
       f.scheduled_arrival,
       timezone(arr.timezone, f.scheduled_arrival) AS scheduled_arrival_local,
       (f.scheduled_arrival - f.scheduled_departure) AS scheduled_duration,
       f.departure_airport,
       dep.airport_name AS departure_airport_name,
       dep.city AS departure_city,
       f.arrival_airport,
       arr.airport_name AS arrival_airport_name,
       arr.city AS arrival_city,
       f.status,
       f.aircraft_code,
       f.actual_departure,
       timezone(dep.timezone, f.actual_departure) AS actual_departure_local,
       f.actual_arrival,
       timezone(arr.timezone, f.actual_arrival) AS actual_arrival_local,
       (f.actual_arrival - f.actual_departure) AS actual_duration
FROM bookings.flights f,
     bookings.airports dep,
     bookings.airports arr
WHERE ((f.departure_airport = dep.airport_code) AND (f.arrival_airport =
arr.airport_code));
```

Now creation of the distributed schema is complete. Let's turn off broadcasting of DDL statements:

```
SET shardman.broadcast_ddl TO off;
```

17.3.2.3.2. Complex Approach — `book_ref` Sharding Key

A more complex approach to the sharding key choice involves the source schema modification, inclusion of new parameters in queries and other important changes.

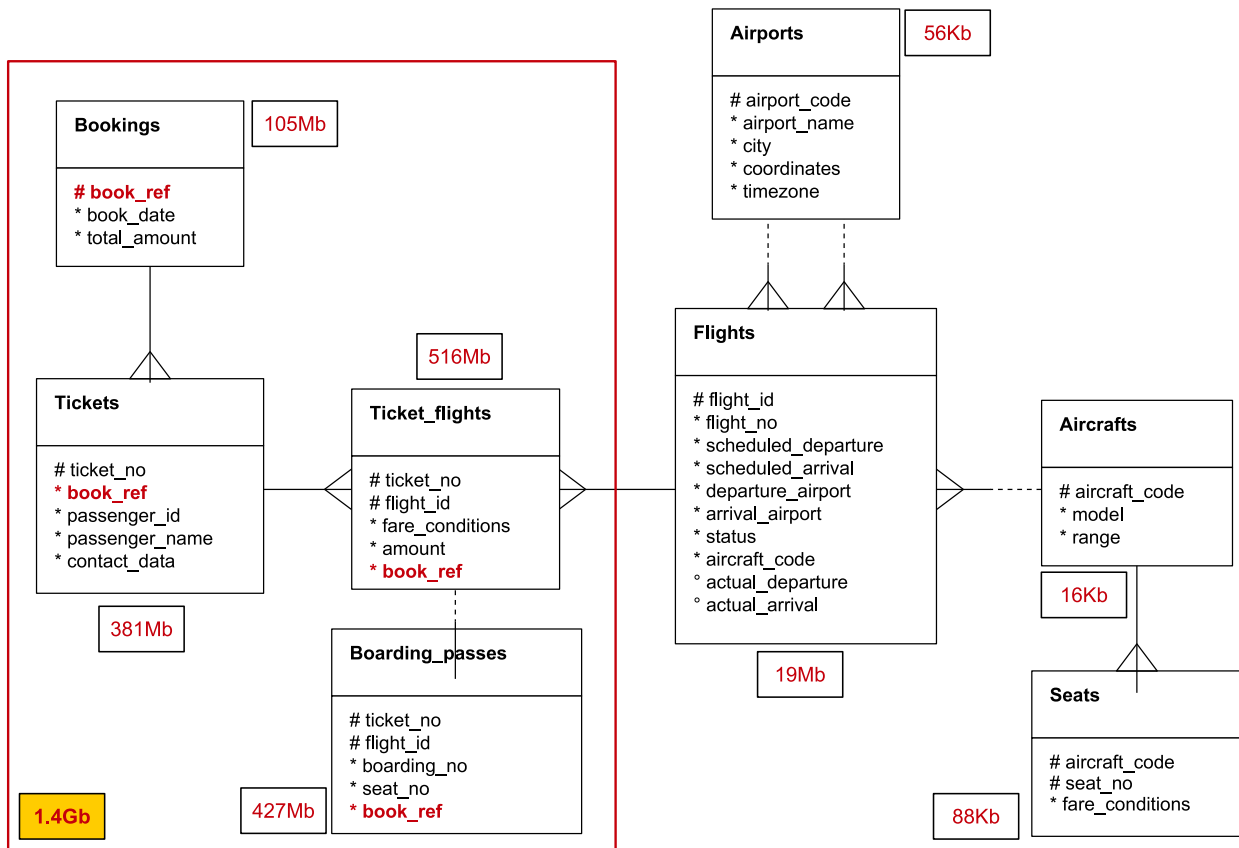
What if an airline is in the market for over 10 years and the `bookings` table reaches the size that does not allow you to continue having it global anymore? But distributing its data is impossible either as it does not contain fields contained in other tables that it can be distributed among (as in [variant 1](#)).

When modifying the source schema, another field can be appropriate for use as a sharding key.

Looking at the `bookings` table, we can notice that values of the `book_ref` field are unique and this field is a primary key. `book_ref` is also a foreign key to the `tickets` table. So this field seems suitable for being the sharding key. However, `book_ref` is missing from the `ticket_flights` and `boarding_passes` tables.

If we add `book_ref` to the `ticket_flights` and `boarding_passes` tables, distributing of all the tables `bookings`, `tickets`, `ticket_flights` and `boarding_passes` with the `book_ref` sharding key becomes possible.

`book_ref` should be added to `ticket_flights` and `boarding_passes` in the source schema, and `book_ref` must be filled with data from the `bookings` table.

Figure 17.4. Source Schema Modification**17.3.2.3.2.1. Modifying the Source Schema**

To properly transfer data from the source schema to the distributed one, the schema should be modified as follows:

1. Add the `book_ref` field to the `ticket_flights` and `boarding_passes` tables:

```
ALTER TABLE ticket_flights
  ADD COLUMN book_ref char(6);
```

```
ALTER TABLE boarding_passes
  ADD COLUMN book_ref char(6);
```

2. In these tables, fill the added `book_ref` field with data:

```
WITH batch AS (SELECT book_ref,
                      ticket_no
                 FROM tickets)
UPDATE ticket_flights
  SET book_ref = batch.book_ref
  FROM batch
 WHERE ticket_flights.ticket_no = batch.ticket_no
    AND ticket_flights.book_ref IS NULL;
```

```
WITH batch AS (SELECT book_ref,
                      ticket_no
                 FROM tickets)
```

```
UPDATE boarding_passes
  SET book_ref = batch.book_ref
  FROM batch
WHERE boarding_passes.ticket_no = batch.ticket_no
  AND boarding_passes.book_ref IS NULL;
```

Avoid using this example in a loaded production system as this approach blocks entire tables, that is, all rows in the tables. In production systems, data should be transferred incrementally, by parts.

Now the database schema is ready for data transferring.

17.3.2.3.2.2. Creating a Schema Distributed by `book_ref`

Here the Postgres Pro Shardman `shardman.broadcast_all_sql` function is used to broadcast DDL statements on all cluster nodes. Let's create the `bookings` schema on all shards:

```
SELECT shardman.broadcast_all_sql('CREATE SCHEMA bookings');
```

As tables in the schema are linked with an external key, the order of creating tables matters.

First we create a utility function `bookings.now()`:

```
SELECT shardman.broadcast_all_sql(
  $sql$
  CREATE FUNCTION bookings.now() RETURNS timestamp with time zone
  LANGUAGE sql IMMUTABLE COST 0.00999999978
  AS
  $q$
  SELECT $qq$2016-10-13 17:00:00$qq$::TIMESTAMP
          AT TIME ZONE $zz$Europe/Moscow$zz$;
  $q$;
  $sql$
);
```

Tables, users and sequences are created with the regular SQL. This function is not needed for that.

In this example, the global sequence is not explicitly created as for the `bigserial` type, Postgres Pro Shardman creates a global sequence automatically.

Now let's create global tables using the following DDL statements:

```
CREATE TABLE bookings.aircrafts (
  aircraft_code character(3) NOT NULL PRIMARY KEY,
  model text NOT NULL,
  range integer NOT NULL,
  CONSTRAINT aircrafts_range_check CHECK ((range > 0))
) WITH (global);

CREATE TABLE bookings.seats (
  aircraft_code character(3) REFERENCES bookings.aircrafts(aircraft_code),
  seat_no character varying(4) NOT NULL,
  fare_conditions character varying(10) NOT NULL,
  CONSTRAINT seats_fare_conditions_check CHECK ((
    (fare_conditions)::text = ANY (ARRAY[
      ('Economy'::character varying)::text,
      ('Comfort'::character varying)::text,
      ('Business'::character varying)::text])
  )),
  PRIMARY KEY (aircraft_code, seat_no)
) WITH (global);
```

```
CREATE TABLE bookings.airports (  
    airport_code character(3) NOT NULL PRIMARY KEY,  
    airport_name text NOT NULL,  
    city text NOT NULL,  
    longitude double precision NOT NULL,  
    latitude double precision NOT NULL,  
    timezone text NOT NULL  
) WITH (global);  
  
CREATE TABLE bookings.flights (  
-- the global sequence will be created automatically  
-- the default value will be assigned  
    flight_id bigserial NOT NULL PRIMARY KEY,  
    flight_no character(6) NOT NULL,  
    scheduled_departure timestamp with time zone NOT NULL,  
    scheduled_arrival timestamp with time zone NOT NULL,  
    departure_airport character(3) REFERENCES bookings.airports(airport_code),  
    arrival_airport character(3) REFERENCES bookings.airports(airport_code),  
    status character varying(20) NOT NULL,  
    aircraft_code character(3) REFERENCES bookings.aircrafts(aircraft_code),  
    actual_departure timestamp with time zone,  
    actual_arrival timestamp with time zone,  
    CONSTRAINT flights_check CHECK ((scheduled_arrival > scheduled_departure)),  
    CONSTRAINT flights_check1 CHECK ((  
        (actual_arrival IS NULL)  
        OR ((actual_departure IS NOT NULL)  
            AND (actual_arrival IS NOT NULL)  
            AND (actual_arrival > actual_departure)))),  
    CONSTRAINT flights_status_check CHECK (  
        ((status)::text = ANY (  
            ARRAY[('On Time'::character varying)::text,  
                ('Delayed'::character varying)::text,  
                ('Departed'::character varying)::text,  
                ('Arrived'::character varying)::text,  
                ('Scheduled'::character varying)::text,  
                ('Cancelled'::character varying)::text])))  
) WITH (global);  
  
ALTER TABLE bookings.flights  
    ADD CONSTRAINT flights_flight_no_scheduled_departure_key  
    UNIQUE (flight_no, scheduled_departure);
```

Now let's create sharded tables `bookings`, `tickets`, `ticket_flights` **and** `boarding_passes` **in the** `bookings` **schema, as in the previous example:**

```
-- no modifications to these tables are done except distributing them  
CREATE TABLE bookings.bookings (  
    book_ref character(6) NOT NULL PRIMARY KEY,  
    book_date timestamp with time zone NOT NULL,  
    total_amount numeric(10,2) NOT NULL  
) WITH (distributed_by='book_ref', num_parts=4);  
  
CREATE TABLE bookings.tickets (  
    ticket_no character(13) NOT NULL,  
    book_ref character(6) REFERENCES bookings.bookings(book_ref),  
    passenger_id character varying(20) NOT NULL,  
    passenger_name text NOT NULL,
```

```

    contact_data jsonb,
    PRIMARY KEY (book_ref, ticket_no)
) WITH (distributed_by='book_ref', colocate_with='bookings.bookings');

-- adding the book_ref foreign key to these tables
CREATE TABLE bookings.ticket_flights (
    ticket_no character(13) NOT NULL,
    flight_id bigint NOT NULL,
    fare_conditions character varying(10) NOT NULL,
    amount numeric(10,2) NOT NULL,
    book_ref character(6) NOT NULL, -- <= added book_ref
    CONSTRAINT ticket_flights_amount_check
        CHECK ((amount >= (0)::numeric)),
    CONSTRAINT ticket_flights_fare_conditions_check
        CHECK (((fare_conditions)::text = ANY (
            ARRAY[('Economy'::character varying)::text,
                ('Comfort'::character varying)::text,
                ('Business'::character varying)::text]))),
    FOREIGN KEY (book_ref, ticket_no)
        REFERENCES bookings.tickets(book_ref, ticket_no),
    PRIMARY KEY (book_ref, ticket_no, flight_id) -- <= changed the primary key
) with (distributed_by='book_ref', colocate_with='bookings.bookings');

CREATE TABLE bookings.boarding_passes (
    ticket_no character(13) NOT NULL,
    flight_id bigint NOT NULL,
    boarding_no integer NOT NULL,
    seat_no character varying(4) NOT NULL,
    book_ref character(6) NOT NULL, -- <= added book_ref
    FOREIGN KEY (book_ref, ticket_no, flight_id)
        REFERENCES bookings.ticket_flights(book_ref, ticket_no, flight_id),
    PRIMARY KEY (book_ref, ticket_no, flight_id)
) WITH (distributed_by='book_ref', colocate_with='bookings.bookings');
-- constraints must contain the sharding key
ALTER TABLE bookings.boarding_passes
    ADD CONSTRAINT boarding_passes_flight_id_boarding_no_key
    UNIQUE (book_ref, ticket_no, flight_id, boarding_no);

ALTER TABLE bookings.boarding_passes
    ADD CONSTRAINT boarding_passes_flight_id_seat_no_key
    UNIQUE (book_ref, ticket_no, flight_id, seat_no);

```

Let's create the bookings.flights view:

```

SELECT shardman.broadcast_all_sql($$
CREATE VIEW bookings.flights_v AS
SELECT f.flight_id,
       f.flight_no,
       f.scheduled_departure,
       timezone(dep.timezone, f.scheduled_departure) AS scheduled_departure_local,
       f.scheduled_arrival,
       timezone(arr.timezone, f.scheduled_arrival)   AS scheduled_arrival_local,
       (f.scheduled_arrival - f.scheduled_departure) AS scheduled_duration,
       f.departure_airport,
       dep.airport_name                               AS departure_airport_name,
       dep.city                                       AS departure_city,
       f.arrival_airport,

```

```
arr.airport_name           AS arrival_airport_name,
arr.city                   AS arrival_city,
f.status,
f.aircraft_code,
f.actual_departure,
timezone(dep.timezone, f.actual_departure) AS actual_departure_local,
f.actual_arrival,
timezone(arr.timezone, f.actual_arrival)   AS actual_arrival_local,
(f.actual_arrival - f.actual_departure)    AS actual_duration
FROM bookings.flights f,
      bookings.airports dep,
      bookings.airports arr
WHERE ((f.departure_airport = dep.airport_code) AND (f.arrival_airport =
arr.airport_code));
$$);
```

The schema creation is now complete. Let's proceed to data migration.

17.3.3. Data Migration

When migrating data, the order of fields in the source and target schema is important. The order and types of fields in the non-distributed and distributed databases must be the same.

The migration utility does exactly what is requested by the user, who does not interfere with data migration processes except, maybe, distributing the data directly to the shard where it must be stored.

Postgres Pro Shardman provides convenient migration tools. Once the distributed schema is created and the sharding key chosen, it is now needed to define the data migration rules. The data source can be either export CSV data files or a single DBMS server.

It is not always convenient to use CSV files as they can reach a pretty large size and require additional resources for storage and transfer.

Migrating data directly from DB to DB without an intermediate storage phase is much more convenient.

The order of loading data during migration must be taken into account. Tables can be linked with a foreign key, so the data in tables that other tables will reference must be loaded first. To follow such an order, in the migration file, you should establish the priority that defines tables whose data must be loaded first. The higher the value of the `priority` parameter, the higher the priority. For example, if the priorities 1, 2 and 3 are defined, tables with the priority 3 will be loaded first, then those with the priority 2, and last with the priority 1.

The `shardmanctl load` command lets you define the order of migrating tables, which can be specified in the configuration YML file.

17.3.3.1. Naive Approach

The following is an example of the `migrate.yml` file:

```
version: "1.0"
migrate:
  connstr: "dbname=demo host=single-pg-instance port=5432 user=postgres
password=*****"
  jobs: 8
  batch: 2000
  options:
  schemas:
    - name: bookings
      # the all parameter set to false turns off automatic creation of pages
      # tables are already created, at the Schema Migration phase
      all: false
```

```
tables:
- name: airports
  # defining a global table
  type: global
  # as tables are linked, data migration priority must be defined
  # setting highest priority to tables whose data
  # must be copied first
  priority: 3
- name: aircrafts
  type: global
  priority: 3
- name: seats
  type: global
  priority: 3
- name: bookings
  type: global
  priority: 3
- name: flights
  type: global
  priority: 3
- name: tickets
  type: sharded
  # defining a sharded table
  # specifying the sharding key
  distributedby: ticket_no
  partitions: 4
  priority: 2
- name: ticket_flights
  type: sharded
  distributedby: ticket_no
  # defining a sharded and colocated table
  # specifying the name of the table that ticket_flights table will be
colocated with
  colocatewith: tickets
  partitions: 4
  priority: 2
- name: boarding_passes
  type: sharded
  distributedby: ticket_no
  colocatewith: tickets
  partitions: 4
  priority: 1
```

This file defines the data source, that is, the `single-pg-instance` node, its connection port, user name and password, and data source DB name. Some parameters of the migration utility operation are also defined (there can be quite a few of them, as explained in [the section called “Loading Data with a Schema from PostgreSQL”](#)). The file also defines the number of threads — 8, batch size, that is, the number of rows organized into batches for processing during migration, as well as table processing priorities. The data for the global tables is migrated first, then the data for the sharded tables `tickets` and `ticket_flights`, and migration of the `boarding_passes` table completes the migration. The value of `priority` defines the priority of data loading, data for tables with higher value will be loaded earlier than with the lower value. The following command performs the migration:

```
shardmanctl load --schema migrate.yml
```

If the utility completes with the message “data loading completed successfully”, it means that the migration was a success.

17.3.3.2. Complex Approach

With this approach, the launch and operation of the [shardmanctl](#) utility in the `load` mode is the same as with the naive approach. However, the file that defines the order of loading tables will slightly differ as the sharding key has changed:

```
---
version: "1.0"
migrate:
  connstr: "dbname=demo host=single-pg-instance port=5432 user=postgres
password=postgres"
  jobs: 8
  batch: 2000
  options:
  schemas:
    - name: bookings
      all: false
      tables:
        - name: airports
          type: global
          priority: 5
        - name: aircrafts
          type: global
          priority: 5
        - name: seats
          type: global
          priority: 5
        - name: flights
          type: global
          priority: 5
        - name: bookings
          type: sharded
          priority: 4
          partitions: 4
          distributedby: book_ref
        - name: tickets
          type: sharded
          distributedby: book_ref
          colocatewith: bookings
          partitions: 4
          priority: 3
        - name: ticket_flights
          type: sharded
          distributedby: book_ref
          colocatewith: bookings
          partitions: 4
          priority: 2
        - name: boarding_passes
          type: sharded
          distributedby: book_ref
          colocatewith: bookings
          partitions: 4
          priority: 1
```

17.3.4. Checking Postgres Pro Shardman Migration

When all the migration operations were performed successfully, it's time to check how queries are executed in the distributed schema.

17.3.4.1. q1 Query

The q1 query is pretty simple, it selects the booking with the specified number:

```
SELECT *
  FROM bookings.bookings b
 WHERE b.book_ref = '0824C5';
```

For the regular PostgreSQL and for the `ticket_no` sharding key, this query runs comparably fast. How fast the query is for the `book_ref` sharding key, depends on the shard where it is executed. If it is executed in a shard where there is physically no data, Postgres Pro Shardman sends the query to another shard, which causes a time delay due to network communication.

17.3.4.2. q2 Query

This q2 query selects all the tickets from the specified booking:

```
SELECT t.*
  FROM bookings.bookings b
 JOIN bookings.tickets t
    ON t.book_ref = b.book_ref
 WHERE b.book_ref = '0824C5';
```

With the `book_ref` sharding key, the query is pushed down to shards and the global table is joined with partitions of a sharded table:

```
Foreign Scan (actual rows=2 loops=1)
  Relations: (bookings_2_fdw b) INNER JOIN (tickets_2_fdw t)
  Network: FDW bytes sent=433 received=237
```

Let's look at the query plan for the `ticket_no` sharding key:

```
Append (actual rows=2 loops=1)
  Network: FDW bytes sent=1263 received=205
  -> Nested Loop (actual rows=1 loops=1)
    -> Seq Scan on tickets_0 t_1 (actual rows=1 loops=1)
        Filter: (book_ref = '0824C5'::bpchar)
        Rows Removed by Filter: 207092
    -> Index Only Scan using bookings_pkey on bookings b (actual rows=1
loops=1)
        Index Cond: (book_ref = '0824C5'::bpchar)
        Heap Fetches: 0
  -> Async Foreign Scan (actual rows=1 loops=1)
        Relations: (tickets_1_fdw t_2) INNER JOIN (bookings b)
        Network: FDW bytes sent=421 received=205
  -> Async Foreign Scan (actual rows=0 loops=1)
        Relations: (tickets_2_fdw t_3) INNER JOIN (bookings b)
        Network: FDW bytes sent=421
  -> Async Foreign Scan (actual rows=0 loops=1)
        Relations: (tickets_3_fdw t_4) INNER JOIN (bookings b)
        Network: FDW bytes sent=421
```

The plan contains `Async Foreign Scan` nodes, which mean network data exchange between the query source node and shards, that is, data is received from shards and final processing is done on the query source node.

Look at the `Network` line. A good criterion of whether query execution on shards is optimal is the value of `received`. The lower its value, the better shards execute distributed queries. Most processing is done remotely, and the query source node gets the result that is ready for further processing.

The case where the sharding key is `book_ref` looks much better as the table with ticket numbers already contains `book_ref`.

The plan of the query to be executed on an arbitrary node is as follows:

```
Foreign Scan (actual rows=2 loops=1)
  Relations: (bookings_2_fdw b) INNER JOIN (tickets_2_fdw t)
  Network: FDW bytes sent=433 received=237
```

The network data exchange is only done with one shard, in which the query is executed. It is `shard-3`, and the `tickets_2` partition of the `tickets` table is on the fourth node.

If this query is executed in the shard where the data is physically located, the query will be executed yet faster.

Let's look at the plan:

```
Nested Loop (actual rows=2 loops=1)
-> Index Only Scan using bookings_2_pkey on bookings_2
-> Bitmap Heap Scan on tickets_2
    -> Bitmap Index Scan on tickets_2_book_ref_idx
```

Network data exchange is not needed here as the requested data is located within the shard in which the query is executed.

In some cases, the choice of the shard for query execution matters. Being aware of the distribution logic, you can implement it at the application level and send some queries immediately to the shard where the needed data is located based on the sharding key.

17.3.4.3. q3 Query

The `q3` query finds all the flights for one of the tickets in the booking selected earlier:

```
SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
  ON tf.ticket_no = t.ticket_no
WHERE t.ticket_no = '0005435126781';
```

To choose a specific shard for query execution, as discussed in [Section 17.3.4.2](#), note that with the `ticket_no` sharding key, the query execution will be more optimal in the shard that contains the partition with the data. The planner knows that the shard contains all the data needed for joining tables, so no network communication between shards will occur.

For the `book_ref` sharding key, note that from the booking number you can compute the ticket number and request it right from the “proper” shard.

So the query is as follows:

```
SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
  ON tf.ticket_no = t.ticket_no
  AND t.book_ref = tf.book_ref
WHERE t.ticket_no = '0005435126781'
AND tf.book_ref = '0824C5';
```

The query is executed more slowly in the shard that does not contain the partition with the data sought:

```
Foreign Scan (actual rows=6 loops=1)
  Relations: (tickets_1_fdw t) INNER JOIN (ticket_flights_1_fdw tf)
  Network: FDW bytes sent=434 received=369
```

Network communication between shards is present in the plan, as it contains the `Foreign Scan` node.

The importance of including the sharding key in a query can be illustrated with the following query for the `book_ref` sharding key:

```
SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
  ON tf.ticket_no = t.ticket_no
WHERE t.ticket_no = '0005435126781'
AND tf.book_ref = '0824C5';
```

Here the sharding key is not included in `join on` purpose. Let's look at the plan:

```
Nested Loop (actual rows=6 loops=1)
  Network: FDW bytes sent=1419 received=600
  -> Foreign Scan on ticket_flights_2_fdw tf (actual rows=6 loops=1)
      Network: FDW bytes sent=381 received=395
  -> Append (actual rows=1 loops=6)
      Network: FDW bytes sent=1038 received=205
      -> Seq Scan on tickets_0 t_1 (actual rows=0 loops=6)
          Filter: (ticket_no = '0005435126781'::bpchar)
          Rows Removed by Filter: 207273
      -> Async Foreign Scan on tickets_1_fdw t_2 (actual rows=0 loops=6)
          Network: FDW bytes sent=346 received=205
      -> Async Foreign Scan on tickets_2_fdw t_3 (actual rows=1 loops=6)
          Network: FDW bytes sent=346
      -> Async Foreign Scan on tickets_3_fdw t_4 (actual rows=0 loops=6)
          Network: FDW bytes sent=346
```

We can notice differences from previous examples. Here the query was executed on all nodes and index was not used, so to return as few as 6 rows, Postgres Pro Shardman had to sequentially scan whole partitions of the `tickets` table, return the result to the query source node and after that perform `join` with the `ticket_flights` table. Async Foreign Scan nodes indicate the sequential scan of the `tickets` table on shards.

17.3.4.4. q4 Query

This query returns all the flights for all the tickets included in a booking. There are several ways to do this: include a subquery in a `WHERE` clause with the booking number, in the `IN` clause, explicitly list ticket numbers or use the `WHERE...OR` clause. Let's check execution of the query for all these variants.

```
SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
  ON tf.ticket_no = t.ticket_no
WHERE t.ticket_no IN (
  SELECT t.ticket_no
  FROM bookings.bookings b
  JOIN bookings.tickets t
    ON t.book_ref = b.book_ref
  WHERE b.book_ref = '0824C5'
);
```

This is just the query from the non-distributed database that we tried to execute. But its execution is equally poor for both sharding keys.

The query plan is like this:

```

Hash Join (actual rows=12 loops=1)
  Hash Cond: (tf.ticket_no = t.ticket_no)
    -> Append (actual rows=2360335 loops=1)
      -> Async Foreign Scan on ticket_flights_0_fdw tf_1 (actual rows=589983
loops=1)
      -> Async Foreign Scan on ticket_flights_1_fdw tf_2 (actual rows=590175
loops=1)
      -> Seq Scan on ticket_flights_2 tf_3 (actual rows=590174 loops=1)
      -> Async Foreign Scan on ticket_flights_3_fdw tf_4 (actual rows=590003
loops=1)
    -> Hash (actual rows=2 loops=1)
      Buckets: 1024 Batches: 1 Memory Usage: 9kB
      -> Hash Semi Join (actual rows=2 loops=1)
        Hash Cond: (t.ticket_no = t_5.ticket_no)
        -> Append (actual rows=829071 loops=1)
          -> Async Foreign Scan on tickets_0_fdw t_1 (actual rows=207273
loops=1)
          -> Async Foreign Scan on tickets_1_fdw t_2 (actual rows=207058
loops=1)
          -> Seq Scan on tickets_2 t_3 (actual rows=207431 loops=1)
          -> Async Foreign Scan on tickets_3_fdw t_4 (actual rows=207309
loops=1)
        -> Hash (actual rows=2 loops=1)
          Buckets: 1024 Batches: 1 Memory Usage: 9kB
          -> Nested Loop (actual rows=2 loops=1)
            -> Index Only Scan using tickets_2_pkey on tickets_2 t_5
            -> Materialize (actual rows=1 loops=2)
              -> Index Only Scan using bookings_2_pkey on
bookings_2 b

```

This plan shows that Postgres Pro Shardman coped with the `WHERE` subquery, then had to request all the rows of the `tickets` and `ticket_flights` tables and then process them on the query source node. This is a really poor performance. Let's try other variants:

For the `ticket_no` sharding key, the query is:

```

SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
  ON tf.ticket_no = t.ticket_no
WHERE t.ticket_no IN ('0005435126781','0005435126782');

```

and the plan is:

```

Append (actual rows=12 loops=1)
  Network: FDW bytes sent=1098 received=1656
  -> Async Foreign Scan (actual rows=6 loops=1)
    Relations: (tickets_0_fdw t_1) INNER JOIN (ticket_flights_0_fdw tf_1)
    Network: FDW bytes sent=549 received=1656
  -> Async Foreign Scan (actual rows=6 loops=1)
    Relations: (tickets_1_fdw t_2) INNER JOIN (ticket_flights_1_fdw tf_2)
    Network: FDW bytes sent=549

```

Everything is pretty good here: the query was executed on two shards of four, and Append of the results received only had to be done.

Let's recall that `book_ref` is contained in both `tickets` and `ticket_flights` tables. So for the `book_ref` sharding key, the query is:

```
SELECT tf.*, t.*
FROM bookings.tickets t
JOIN bookings.ticket_flights tf
ON tf.ticket_no = t.ticket_no
AND tf.book_ref = t.book_ref
WHERE t.book_ref = '0824C5';
```

and the plan is:

```
Foreign Scan (actual rows=12 loops=1)
  Relations: (tickets_2_fdw t) INNER JOIN (ticket_flights_2_fdw tf)
  Network: FDW bytes sent=547 received=1717
```

This is an excellent result — the query was modified to execute well in the distributed schema.

17.3.4.5. q5 Query

This is a small analytical query, which returns the names and ticket numbers of the passengers who got registered first.

```
SELECT t.passenger_name, t.ticket_no
FROM bookings.tickets t
JOIN bookings.boarding_passes bp
  ON bp.ticket_no = t.ticket_no
GROUP BY t.passenger_name, t.ticket_no
HAVING max(bp.boarding_no) = 1
AND count(*) > 1;
```

This query is executed pretty slowly for both sharding keys. Below is the plan for `book_ref`:

```
HashAggregate (actual rows=424 loops=1)
  Group Key: t.ticket_no
  Filter: ((max(bp.boarding_no) = 1) AND (count(*) > 1))
  Batches: 85  Memory Usage: 4265kB  Disk Usage: 112008kB
  Rows Removed by Filter: 700748
  Network: FDW bytes sent=1215 received=77111136
  -> Append (actual rows=1894295 loops=1)
    Network: FDW bytes sent=1215 received=77111136
    -> Async Foreign Scan (actual rows=473327 loops=1)
      Relations: (tickets_0_fdw t_1) INNER JOIN (boarding_passes_0_fdw
bp_1)
      Network: FDW bytes sent=404 received=813128
    -> Async Foreign Scan (actual rows=472632 loops=1)
      Relations: (tickets_1_fdw t_2) INNER JOIN (boarding_passes_1_fdw
bp_2)
      Network: FDW bytes sent=404
    -> Async Foreign Scan (actual rows=475755 loops=1)
      Relations: (tickets_2_fdw t_3) INNER JOIN (boarding_passes_2_fdw
bp_3)
      Network: FDW bytes sent=407
    -> Hash Join (actual rows=472581 loops=1)
      Hash Cond: (bp_4.ticket_no = t_4.ticket_no)
      Network: FDW bytes received=28841344
    -> Seq Scan on boarding_passes_3 bp_4 (actual rows=472581 loops=1)
    -> Hash (actual rows=207118 loops=1)
      Buckets: 65536  Batches: 4  Memory Usage: 3654kB
      Network: FDW bytes received=9176680
    -> Seq Scan on tickets_3 t_4 (actual rows=207118 loops=1)
      Network: FDW bytes received=9176680
```

Note a pretty large amount of network data transfer between shards. Let's improve the query by adding `book_ref` as one more condition for joining tables:

```
SELECT t.passenger_name, t.ticket_no
FROM bookings.tickets t
JOIN bookings.boarding_passes bp
  ON bp.ticket_no = t.ticket_no
  AND bp.book_ref=t.book_ref -- <= added book_ref
GROUP BY t.passenger_name, t.ticket_no
HAVING max(bp.boarding_no) = 1
AND count(*) > 1;
```

Let's look at the query plan:

```
GroupAggregate (actual rows=424 loops=1)
  Group Key: t.passenger_name, t.ticket_no
  Filter: ((max(bp.boarding_no) = 1) AND (count(*) > 1))
  Rows Removed by Filter: 700748
  Network: FDW bytes sent=1424 received=77092816
  -> Merge Append (actual rows=1894295 loops=1)
    Sort Key: t.passenger_name, t.ticket_no
    Network: FDW bytes sent=1424 received=77092816
    -> Foreign Scan (actual rows=472757 loops=1)
      Relations: (tickets_0_fdw t_1) INNER JOIN (boarding_passes_0_fdw
bp_1)
      Network: FDW bytes sent=472 received=2884064
    -> Sort (actual rows=472843 loops=1)
      Sort Key: t_2.passenger_name, t_2.ticket_no
      Sort Method: external merge  Disk: 21152kB
      Network: FDW bytes received=22753536
    -> Hash Join (actual rows=472843 loops=1)
      Hash Cond: ((bp_2.ticket_no = t_2.ticket_no) AND (bp_2.book_ref
= t_2.book_ref))
      Network: FDW bytes received=22753536
    -> Seq Scan on boarding_passes_1 bp_2 (actual rows=472843
loops=1)
    -> Hash (actual rows=207058 loops=1)
      Buckets: 65536  Batches: 8  Memory Usage: 2264kB
      Network: FDW bytes received=22753536
    -> Seq Scan on tickets_1 t_2 (actual rows=207058
loops=1)
      Network: FDW bytes received=22753536
    -> Foreign Scan (actual rows=474715 loops=1)
      Relations: (tickets_2_fdw t_3) INNER JOIN (boarding_passes_2_fdw
bp_3)
      Network: FDW bytes sent=476 received=2884120
    -> Foreign Scan (actual rows=473980 loops=1)
      Relations: (tickets_3_fdw t_4) INNER JOIN (boarding_passes_3_fdw
bp_4)
      Network: FDW bytes sent=476 received=25745384
```

The situation considerably improved, the result was received on the query source node, and then final filtering, grouping and joining data were done.

For the `ticket_no` sharding key, the source query plan looks like this:

```
HashAggregate (actual rows=424 loops=1)
```

```

Group Key: t.ticket_no
Filter: ((max(bp.boarding_no) = 1) AND (count(*) > 1))
Batches: 85  Memory Usage: 4265kB  Disk Usage: 111824kB
Rows Removed by Filter: 700748
Network: FDW bytes sent=1188 received=77103620
-> Append (actual rows=1894295 loops=1)
    Network: FDW bytes sent=1188 received=77103620
    -> Async Foreign Scan (actual rows=473327 loops=1)
        Relations: (tickets_0_fdw t_1) INNER JOIN (boarding_passes_0_fdw
bp_1)
    Network: FDW bytes sent=394
    -> Hash Join (actual rows=472632 loops=1)
        Hash Cond: (bp_2.ticket_no = t_2.ticket_no)
        Network: FDW bytes received=77103620
        -> Seq Scan on boarding_passes_1 bp_2 (actual rows=472632 loops=1)
        -> Hash (actual rows=206712 loops=1)
            Buckets: 65536  Batches: 4  Memory Usage: 3654kB
            Network: FDW bytes received=23859576
            -> Seq Scan on tickets_1 t_2 (actual rows=206712 loops=1)
                Network: FDW bytes received=23859576
        -> Async Foreign Scan (actual rows=475755 loops=1)
            Relations: (tickets_2_fdw t_3) INNER JOIN (boarding_passes_2_fdw
bp_3)
    Network: FDW bytes sent=397
    -> Async Foreign Scan (actual rows=472581 loops=1)
        Relations: (tickets_3_fdw t_4) INNER JOIN (boarding_passes_3_fdw
bp_4)
    Network: FDW bytes sent=397

```

We can see that table joining is done on shards, while data filtering, grouping and aggregation are done on the query source node. The source query does not need to be modified in this case.

17.3.4.6. q6 Query

For each ticket booked a week ago from now, this query displays all the included flight segments, together with connection time.

```

SELECT tf.ticket_no,f.departure_airport,
       f.arrival_airport,f.scheduled_arrival,
       lead(f.scheduled_departure) OVER w AS next_departure,
       lead(f.scheduled_departure) OVER w - f.scheduled_arrival AS gap
FROM bookings.bookings b
JOIN bookings.tickets t
ON t.book_ref = b.book_ref
JOIN bookings.ticket_flights tf
ON tf.ticket_no = t.ticket_no
JOIN bookings.flights f
ON tf.flight_id = f.flight_id
WHERE b.book_date = bookings.now()::date - INTERVAL '7 day'

WINDOW w AS (
PARTITION BY tf.ticket_no
ORDER BY f.scheduled_departure);

```

For this query, the type of the `book_date` column must be cast from the `timestamp_tz` to `date`. When casting types, PostgreSQL casts the column data type to the data type specified in the filtering condition, but not vice versa. Therefore, Postgres Pro Shardman must first get all the data from other shards, cast the type and apply filtering only after that. The query plan looks like this:

```

WindowAgg (actual rows=26 loops=1)
  Network: FDW bytes sent=1750 received=113339240
  -> Sort (actual rows=26 loops=1)
    Sort Key: tf.ticket_no, f.scheduled_departure
    Sort Method: quicksort  Memory: 27kB
    Network: FDW bytes sent=1750 received=113339240
    -> Append (actual rows=26 loops=1)
      Network: FDW bytes sent=1750 received=113339240
      -> Hash Join (actual rows=10 loops=1)
        Hash Cond: (t_1.book_ref = b.book_ref)
        Network: FDW bytes sent=582 received=37717376
      -> Hash Join (actual rows=6 loops=1)
        Hash Cond: (t_2.book_ref = b.book_ref)
        Network: FDW bytes sent=582 received=37700608
      -> Hash Join (actual rows=2 loops=1)
        Hash Cond: (t_3.book_ref = b.book_ref)
        Network: FDW bytes sent=586 received=37921256
      -> Nested Loop (actual rows=8 loops=1)
        -> Nested Loop (actual rows=8 loops=1)
          -> Hash Join (actual rows=2 loops=1)
            Hash Cond: (t_4.book_ref = b.book_ref)
          -> Seq Scan on tickets_3 t_4 (actual rows=207118
loops=1)
          -> Index Scan using flights_pkey on flights f (actual rows=1
loops=8)
            Index Cond: (flight_id = tf_4.flight_id)

```

Pay attention to the number of bytes received from other cluster shards and to the sequential scan of the `tickets` table. Let's try to rewrite the query to avoid the type cast.

The idea is pretty simple: the interval will be computed at the application level rather than at the database level, and the data of the `timestampz` type will be readily passed to the query. Besides, creation of an additional index can help:

```
CREATE INDEX if not exists bookings_date_idx ON bookings.bookings(book_date);
```

For the `book_ref` sharding key, the query looks like this:

```

SELECT tf.ticket_no, f.departure_airport,
       f.arrival_airport, f.scheduled_arrival,
       lead(f.scheduled_departure) OVER w AS next_departure,
       lead(f.scheduled_departure) OVER w - f.scheduled_arrival AS gap
FROM bookings.bookings b
JOIN bookings.tickets t
ON t.book_ref = b.book_ref
JOIN bookings.ticket_flights tf
ON tf.ticket_no = t.ticket_no
AND tf.book_ref = t.book_ref -- <= added book_ref
JOIN bookings.flights f
ON tf.flight_id = f.flight_id
WHERE b.book_date = '2016-10-06 14:00:00+00'
WINDOW w AS (
PARTITION BY tf.ticket_no
ORDER BY f.scheduled_departure);

```

This query has a different plan:

```
WindowAgg (actual rows=18 loops=1)
```



```

Network: FDW bytes sent=2268 received=892
-> Sort (actual rows=18 loops=1)
    Sort Key: tf.ticket_no, f.scheduled_departure
    Sort Method: quicksort  Memory: 26kB
    Network: FDW bytes sent=2268 received=892
-> Append (actual rows=18 loops=1)
    Network: FDW bytes sent=2268 received=892
    -> Nested Loop (actual rows=4 loops=1)
        -> Nested Loop (actual rows=4 loops=1)
            -> Nested Loop (actual rows=1 loops=1)
                -> Bitmap Heap Scan on bookings_0 b_1
                    Heap Blocks: exact=1
                -> Bitmap Index Scan on
bookings_0_book_date_idx
                                -> Index Only Scan using tickets_0_pkey on
tickets_0 t_1
                                    Index Cond: (book_ref = b_1.book_ref)
                                    Heap Fetches: 0
                                -> Index Only Scan using ticket_flights_0_pkey on
ticket_flights_0 tf_1
                                    Heap Fetches: 0
                                -> Index Scan using flights_pkey on flights f (actual rows=1
loops=4)
                                    Index Cond: (flight_id = tf_1.flight_id)
                                -> Async Foreign Scan (actual rows=14 loops=1)
                                    Network: FDW bytes sent=754 received=892
                                -> Async Foreign Scan (actual rows=0 loops=1)
                                    Network: FDW bytes sent=757 -- received=0!
                                -> Async Foreign Scan (actual rows=0 loops=1)
                                    Network: FDW bytes sent=757 -- received=0!

```

This is much better. First, the whole table is not scanned, Index Only Scan is only included. Second, it is clear how much the amount of network data transfer between nodes is reduced.

17.3.4.7. q7 Query

Assume that statistics is needed showing how many passengers there are per booking. To find this out, let's first compute the number of passengers in each booking and then the number of bookings with each number of passengers.

```

SELECT tt.cnt, count(*)
FROM (
    SELECT count(*) cnt
    FROM bookings.tickets t
    GROUP BY t.book_ref
) tt
GROUP BY tt.cnt
ORDER BY tt.cnt;

```

This query processes all the data in the `tickets` and `bookings` tables. So intensive network data exchange between shards cannot be avoided. Also note that the value of the `work_mem` parameter must be pretty high to avoid the use of disk when joining tables. So let's change the value of `work_mem` in the cluster:

```
shardmanctl set work_mem='256MB';
```

The query plan for the `ticket_no` sharding key is as follows:

```

GroupAggregate (actual rows=5 loops=1)
  Group Key: tt.cnt
  Network: FDW bytes sent=798 received=18338112
  -> Sort (actual rows=593433 loops=1)
    Sort Key: tt.cnt
    Sort Method: quicksort  Memory: 57030kB
    Network: FDW bytes sent=798 received=18338112
    -> Subquery Scan on tt (actual rows=593433 loops=1)
      Network: FDW bytes sent=798 received=18338112
      -> Finalize HashAggregate (actual rows=593433 loops=1)
        Group Key: t.book_ref
        Batches: 1  Memory Usage: 81953kB
        Network: FDW bytes sent=798 received=18338112
        -> Append (actual rows=763806 loops=1)
          Network: FDW bytes sent=798 received=18338112
          -> Async Foreign Scan (actual rows=190886 loops=1)
            Relations: Aggregate on (tickets_0_fdw t)
            Network: FDW bytes sent=266 received=1558336
          -> Async Foreign Scan (actual rows=190501 loops=1)
            Relations: Aggregate on (tickets_1_fdw t_1)
            Network: FDW bytes sent=266
          -> Async Foreign Scan (actual rows=191589 loops=1)
            Relations: Aggregate on (tickets_2_fdw t_2)
            Network: FDW bytes sent=266
          -> Partial HashAggregate (actual rows=190830 loops=1)
            Group Key: t_3.book_ref
            Batches: 1  Memory Usage: 36881kB
            Network: FDW bytes received=4981496
            -> Seq Scan on tickets_3 t_3 (actual rows=207118
loops=1)
                                     Network: FDW bytes received=4981496

```

The query plan for the `book_ref` sharding key is as follows:

```

Sort (actual rows=5 loops=1)
  Sort Key: (count(*))
  Sort Method: quicksort  Memory: 25kB
  Network: FDW bytes sent=798 received=14239951
  -> HashAggregate (actual rows=5 loops=1)
    Group Key: (count(*))
    Batches: 1  Memory Usage: 40kB
    Network: FDW bytes sent=798 received=14239951
    -> Append (actual rows=593433 loops=1)
      Network: FDW bytes sent=798 received=14239951
      -> GroupAggregate (actual rows=148504 loops=1)
        Group Key: t.book_ref
        -> Index Only Scan using tickets_0_book_ref_idx on tickets_0 t
(rows=207273)
                                     Heap Fetches: 0
      -> Async Foreign Scan (actual rows=148256 loops=1)
        Relations: Aggregate on (tickets_1_fdw t_1)
        Network: FDW bytes sent=266 received=1917350
      -> Async Foreign Scan (actual rows=148270 loops=1)
        Relations: Aggregate on (tickets_2_fdw t_2)
        Network: FDW bytes sent=266
      -> Async Foreign Scan (actual rows=148403 loops=1)
        Relations: Aggregate on (tickets_3_fdw t_3)
        Network: FDW bytes sent=266

```

The query plans differ first by the order of joining tables and by the computation of aggregates.

For the `ticket_no` sharding key, all the partially aggregated data of the joined tables is received (17 Mb), and all the rest of processing is performed on the query source node.

For the `book_ref` sharding key, as it is included in the query, most of the computation of aggregates is performed on the nodes and only the result (13 Mb) is returned to the query source node, which is then finalized.

17.3.4.8. q8 Query

This query answers the question: which are the most frequent combinations of first and last names in bookings and what is the ratio of the passengers with such names to the total number of passengers. A window function is used to get the result:

```
SELECT passenger_name,
       round( 100.0 * cnt / sum(cnt) OVER (), 2)
       AS percent
FROM (
  SELECT passenger_name,
         count(*) cnt
  FROM bookings.tickets
  GROUP BY passenger_name
) t
ORDER BY percent DESC;
```

For both sharding keys, the query plan looks like this:

```
Sort (actual rows=27909 loops=1)
  Sort Key: (round(((100.0 * ((count(*)))::numeric) / sum((count(*))) OVER (?)),
2)) DESC
  Sort Method: quicksort  Memory: 3076kB
  Network: FDW bytes sent=816 received=2376448
  -> WindowAgg (actual rows=27909 loops=1)
    Network: FDW bytes sent=816 received=2376448
    -> Finalize HashAggregate (actual rows=27909 loops=1)
      Group Key: tickets.passenger_name
      Batches: 1  Memory Usage: 5649kB
      Network: FDW bytes sent=816 received=2376448
      -> Append (actual rows=74104 loops=1)
        Network: FDW bytes sent=816 received=2376448
        -> Partial HashAggregate (actual rows=18589 loops=1)
          Group Key: tickets.passenger_name
          Batches: 1  Memory Usage: 2833kB
          -> Seq Scan on tickets_0 tickets (actual rows=207273
loops=1)

        -> Async Foreign Scan (actual rows=18435 loops=1)
          Relations: Aggregate on (tickets_1_fdw tickets_1)
          Network: FDW bytes sent=272 received=2376448
        -> Async Foreign Scan (actual rows=18567 loops=1)
          Relations: Aggregate on (tickets_2_fdw tickets_2)
          Network: FDW bytes sent=272
        -> Async Foreign Scan (actual rows=18513 loops=1)
          Relations: Aggregate on (tickets_3_fdw tickets_3)
          Network: FDW bytes sent=272
```

The plan shows that the data preprocessing, table joins and partial aggregation are performed on shards, while the final processing is performed on the query source node.

17.3.4.9. q9 Query

This query answers the question: who traveled from Moscow (SVO) to Novosibirsk (OVB) on seat 1A the day before yesterday, and when was the ticket booked. The day before yesterday is computed from the function `booking.now` rather than from the current date. The query in the non-distributed schema is as follows:

```
SELECT
  t.passenger_name,
  b.book_date v
FROM bookings b
JOIN tickets t ON
  t.book_ref = b.book_ref
JOIN boarding_passes bp
  ON bp.ticket_no = t.ticket_no
JOIN flights f ON
  f.flight_id = bp.flight_id
WHERE f.departure_airport = 'SVO'
AND f.arrival_airport = 'OVB'
AND f.scheduled_departure::date = bookings.now()::date - INTERVAL '2 day'
AND bp.seat_no = '1A';
```

As explained for the [q6 Query](#), `INTERVAL` causes the type cast. Let's get rid of it and rewrite the query for the `book_ref` sharding key as follows:

```
SELECT
  t.passenger_name,
  b.book_date v
FROM bookings b
JOIN tickets t ON
  t.book_ref = b.book_ref
JOIN boarding_passes bp
  ON bp.ticket_no = t.ticket_no
  AND bp.book_ref = b.book_ref -- <= added book_ref
JOIN flights f ON
  f.flight_id = bp.flight_id
WHERE f.departure_airport = 'SVO'
AND f.arrival_airport = 'OVB'
AND f.scheduled_departure
  BETWEEN '2016-10-11 14:00:00+00' AND '2016-10-13 14:00:00+00'
AND bp.seat_no = '1A';
```

Let's also create a couple of additional indexes:

```
CREATE INDEX idx_boarding_passes_seats
  ON boarding_passes((seat_no::text));
CREATE INDEX idx_flights_sched_dep
  ON flights(departure_airport, arrival_airport, scheduled_departure);
```

As a result, the query plan appears pretty good:

```
Append (actual rows=1 loops=1)
  Network: FDW bytes sent=2484 received=102
-> Nested Loop (actual rows=1 loops=1)
    Join Filter: (bp_1.ticket_no = t_1.ticket_no)
    Rows Removed by Join Filter: 1
-> Nested Loop (actual rows=1 loops=1)
    -> Hash Join (actual rows=1 loops=1)
        Hash Cond: (bp_1.flight_id = f.flight_id)
```

```

-> Bitmap Heap Scan on boarding_passes_0 bp_1 (actual
rows=4919 loops=1)
    Recheck Cond: ((seat_no)::text = '1A'::text)
    Heap Blocks: exact=2632
    -> Bitmap Index Scan on boarding_passes_0_seat_no_idx
(actual rows=4919)
        Index Cond: ((seat_no)::text = '1A'::text)
-> Hash (actual rows=2 loops=1)
    Buckets: 1024 Batches: 1 Memory Usage: 9kB
    -> Bitmap Heap Scan on flights f (actual rows=2 loops=1)
        Recheck Cond:
        ((departure_airport = 'SVO'::bpchar) AND (arrival_airport =
'OVV'::bpchar) AND
        (scheduled_departure >= '2016-10-11 14:00:00+00'::timestamp
with time zone) AND
        (scheduled_departure < '2016-10-13 14:00:00+00'::timestamp with
time zone))
        Heap Blocks: exact=2
        -> Bitmap Index Scan on idx_flights_sched_dep
(actual rows=2 loops=1)
            Index Cond:
            ((departure_airport = 'SVO'::bpchar) AND
            (arrival_airport = 'OVV'::bpchar) AND
            (scheduled_departure >= '2016-10-11 14:00:00+00'::timestamp
with time zone) AND
            (scheduled_departure <= '2016-10-13 14:00:00+00'::timestamp
with time zone))
-> Index Scan using bookings_0_pkey on bookings_0 b_1 (actual rows=1
loops=1)
    Index Cond: (book_ref = bp_1.book_ref)
-> Index Scan using tickets_0_book_ref_idx on tickets_0 t_1 (actual rows=2
loops=1)
    Index Cond: (book_ref = b_1.book_ref)
-> Async Foreign Scan (actual rows=0 loops=1)
    Relations: (((boarding_passes_1_fdw bp_2) INNER JOIN (flights f)) INNER
JOIN (tickets_1_fdw t_2)) INNER JOIN (bookings_1_fdw b_2)
    Network: FDW bytes sent=826 received=68
-> Async Foreign Scan (actual rows=0 loops=1)
    Relations: (((boarding_passes_2_fdw bp_3) INNER JOIN (flights f)) INNER
JOIN (tickets_2_fdw t_3)) INNER JOIN (bookings_2_fdw b_3)
    Network: FDW bytes sent=829 received=34
-> Async Foreign Scan (actual rows=0 loops=1)
    Relations: (((boarding_passes_3_fdw bp_4) INNER JOIN (flights f)) INNER
JOIN (tickets_3_fdw t_4)) INNER JOIN (bookings_3_fdw b_4)
    Network: FDW bytes sent=829

```

It is clear from this plan that all the table joining was done on shards and the query source node received the result that did not contain rows as the data was located on one shard where the query was executed.

If this query were executed on a different shard, the plan would be the same, but the data for finalization would be received from the shard with the data.

Chapter 18. Server Setup and Operation

This chapter discusses how to set up and run the database server, and its interactions with the operating system.

The directions in this chapter assume that you are working with plain PostgreSQL without any additional infrastructure, for example a copy that you built from source according to the directions in the preceding chapters. If you are working with a pre-packaged or vendor-supplied version of PostgreSQL, it is likely that the packager has made special provisions for installing and starting the database server according to your system's conventions. Consult the package-level documentation for details.

18.1. The Postgres Pro User Account

As with any server daemon that is accessible to the outside world, it is advisable to run Postgres Pro under a separate user account. This user account should only own the data that is managed by the server, and should not be shared with other daemons. (For example, using the user `nobody` is a bad idea.) In particular, it is advisable that this user account not own the Postgres Pro executable files, to ensure that a compromised server process could not modify those executables.

Pre-packaged versions of PostgreSQL will typically create a suitable user account automatically during package installation.

To add a Unix user account to your system, look for a command `useradd` or `adduser`. The user name `postgres` is often used, and is assumed throughout this book, but you can use another name if you like.

18.2. Creating a Database Cluster

For distributed system specifics, refer to the [shardmanctl](#) tool.

Before you can do anything, you must initialize a database storage area on disk. We call this a *database cluster*. (The SQL standard uses the term *catalog cluster*.) A database cluster is a collection of databases that is managed by a single instance of a running database server. After initialization, a database cluster will contain a database named `postgres`, which is meant as a default database for use by utilities, users and third party applications. The database server itself does not require the `postgres` database to exist, but many external utility programs assume it exists. There are two more databases created within each cluster during initialization, named `template1` and `template0`. As the names suggest, these will be used as templates for subsequently-created databases; they should not be used for actual work. (See [Chapter 22](#) for information about creating new databases within a cluster.)

In file system terms, a database cluster is a single directory under which all data will be stored. We call this the *data directory* or *data area*. It is completely up to you where you choose to store your data. There is no default, although locations such as `/usr/local/pgsql/data` or `/var/lib/pgsql/data` are popular. The data directory must be initialized before being used, using the program [initdb](#) which is installed with Postgres Pro.

If you are using a pre-packaged version of PostgreSQL, it may well have a specific convention for where to place the data directory, and it may also provide a script for creating the data directory. In that case you should use that script in preference to running `initdb` directly. Consult the package-level documentation for details.

To initialize a database cluster manually, run `initdb` and specify the desired file system location of the database cluster with the `-D` option, for example:

```
$ initdb -D /usr/local/pgsql/data
```

Note that you must execute this command while logged into the Postgres Pro user account, which is described in the previous section.

Tip

As an alternative to the `-D` option, you can set the environment variable `PGDATA`.

Alternatively, you can run `initdb` via the `pg_ctl` program like so:

```
$ pg_ctl -D /usr/local/pgsql/data initdb
```

This may be more intuitive if you are using `pg_ctl` for starting and stopping the server (see [Section 18.3](#)), so that `pg_ctl` would be the sole command you use for managing the database server instance.

`initdb` will attempt to create the directory you specify if it does not already exist. Of course, this will fail if `initdb` does not have permissions to write in the parent directory. It's generally recommendable that the Postgres Pro user own not just the data directory but its parent directory as well, so that this should not be a problem. If the desired parent directory doesn't exist either, you will need to create it first, using root privileges if the grandparent directory isn't writable. So the process might look like this:

```
root# mkdir /usr/local/pgsql
root# chown postgres /usr/local/pgsql
root# su postgres
postgres$ initdb -D /usr/local/pgsql/data
```

`initdb` will refuse to run if the data directory exists and already contains files; this is to prevent accidentally overwriting an existing installation.

Because the data directory contains all the data stored in the database, it is essential that it be secured from unauthorized access. `initdb` therefore revokes access permissions from everyone but the Postgres Pro user, and optionally, group. Group access, when enabled, is read-only. This allows an unprivileged user in the same group as the cluster owner to take a backup of the cluster data or perform other operations that only require read access.

Note that enabling or disabling group access on an existing cluster requires the cluster to be shut down and the appropriate mode to be set on all directories and files before restarting Postgres Pro. Otherwise, a mix of modes might exist in the data directory. For clusters that allow access only by the owner, the appropriate modes are `0700` for directories and `0600` for files. For clusters that also allow reads by the group, the appropriate modes are `0750` for directories and `0640` for files.

However, while the directory contents are secure, the default client authentication setup allows any local user to connect to the database and even become the database superuser. If you do not trust other local users, we recommend you use one of `initdb`'s `-W`, `--pwprompt` or `--pwfile` options to assign a password to the database superuser. Also, specify `-A scram-sha-256` so that the default `trust` authentication mode is not used; or modify the generated `pg_hba.conf` file after running `initdb`, but *before* you start the server for the first time. (Other reasonable approaches include using `peer` authentication or file system permissions to restrict connections. See [Chapter 20](#) for more information.)

`initdb` also initializes the default locale for the database cluster. Normally, it will just take the locale settings in the environment and apply them to the initialized database. It is possible to specify a different locale for the database; more information about that can be found in [Section 23.1](#). The default sort order used within the particular database cluster is set by `initdb`, and while you can create new databases using different sort order, the order used in the template databases that `initdb` creates cannot be changed without dropping and recreating them. There is also a performance impact for using locales other than `C` or `POSIX`. Therefore, it is important to make this choice correctly the first time.

`initdb` also sets the default character set encoding for the database cluster. Normally this should be chosen to match the locale setting. For details see [Section 23.3](#).

Non-`C` and non-`POSIX` locales rely on the operating system's collation library for character set ordering. This controls the ordering of keys stored in indexes. For this reason, a cluster cannot switch to an

incompatible collation library version, either through snapshot restore, binary streaming replication, a different operating system, or an operating system upgrade.

18.2.1. Use of Secondary File Systems

Many installations create their database clusters on file systems (volumes) other than the machine's "root" volume. If you choose to do this, it is not advisable to try to use the secondary volume's topmost directory (mount point) as the data directory. Best practice is to create a directory within the mount-point directory that is owned by the Postgres Pro user, and then create the data directory within that. This avoids permissions problems, particularly for operations such as `pg_upgrade`, and it also ensures clean failures if the secondary volume is taken offline.

18.2.2. File Systems

Generally, any file system with POSIX semantics can be used for Postgres Pro. Users prefer different file systems for a variety of reasons, including vendor support, performance, and familiarity. Experience suggests that, all other things being equal, one should not expect major performance or behavior changes merely from switching file systems or making minor file system configuration changes.

18.2.2.1. NFS

It is possible to use an NFS file system for storing the Postgres Pro data directory. Postgres Pro does nothing special for NFS file systems, meaning it assumes NFS behaves exactly like locally-connected drives. Postgres Pro does not use any functionality that is known to have nonstandard behavior on NFS, such as file locking.

The only firm requirement for using NFS with Postgres Pro is that the file system is mounted using the `hard` option. With the `hard` option, processes can "hang" indefinitely if there are network problems, so this configuration will require a careful monitoring setup. The `soft` option will interrupt system calls in case of network problems, but Postgres Pro will not repeat system calls interrupted in this way, so any such interruption will result in an I/O error being reported.

It is not necessary to use the `sync` mount option. The behavior of the `async` option is sufficient, since Postgres Pro issues `fsync` calls at appropriate times to flush the write caches. (This is analogous to how it works on a local file system.) However, it is strongly recommended to use the `sync` export option on the NFS *server* on systems where it exists (mainly Linux). Otherwise, an `fsync` or equivalent on the NFS client is not actually guaranteed to reach permanent storage on the server, which could cause corruption similar to running with the parameter `fsync` off. The defaults of these mount and export options differ between vendors and versions, so it is recommended to check and perhaps specify them explicitly in any case to avoid any ambiguity.

In some cases, an external storage product can be accessed either via NFS or a lower-level protocol such as iSCSI. In the latter case, the storage appears as a block device and any available file system can be created on it. That approach might relieve the DBA from having to deal with some of the idiosyncrasies of NFS, but of course the complexity of managing remote storage then happens at other levels.

18.3. Starting the Database Server

Important

In binary installations on Linux systems, the default database is located in `/var/lib/pg-pro/sdm-17/data`, unless you specify a custom directory. See [Section 17.1](#) for details.

For distributed system specifics, refer to the [shardmanctl](#) tool.

Before anyone can access the database, you must start the database server. The database server program is called `postgres`.

If you are using a pre-packaged version of PostgreSQL, it almost certainly includes provisions for running the server as a background task according to the conventions of your operating system. Using the package's infrastructure to start the server will be much less work than figuring out how to do this yourself. Consult the package-level documentation for details.

The bare-bones way to start the server manually is just to invoke `postgres` directly, specifying the location of the data directory with the `-D` option, for example:

```
$ postgres -D /usr/local/pgsql/data
```

which will leave the server running in the foreground. This must be done while logged into the Postgres Pro user account. Without `-D`, the server will try to use the data directory named by the environment variable `PGDATA`. If that variable is not provided either, it will fail.

Normally it is better to start `postgres` in the background. For this, use the usual Unix shell syntax:

```
$ postgres -D /usr/local/pgsql/data >logfile 2>&1 &
```

It is important to store the server's stdout and stderr output somewhere, as shown above. It will help for auditing purposes and to diagnose problems. (See [Section 24.4](#) for a more thorough discussion of log file handling.)

The `postgres` program also takes a number of other command-line options. For more information, see the [postgres](#) reference page and [Chapter 19](#) below.

This shell syntax can get tedious quickly. Therefore the wrapper program `pg_ctl` is provided to simplify some tasks. For example:

```
pg_ctl start -l logfile
```

will start the server in the background and put the output into the named log file. The `-D` option has the same meaning here as for `postgres`. `pg_ctl` is also capable of stopping the server.

Normally, you will want to start the database server when the computer boots. Autostart scripts are operating-system-specific. There are a few example scripts distributed with PostgreSQL in the `contrib/start-scripts` directory. Installing one will require root privileges.

Different systems have different conventions for starting up daemons at boot time. Many systems have a file `/etc/rc.local` or `/etc/rc.d/rc.local`. Others use `init.d` or `rc.d` directories. Whatever you do, the server must be run by the Postgres Pro user account *and not by root* or any other user. Therefore you probably should form your commands using `su postgres -c '...'`. For example:

```
su postgres -c 'pg_ctl start -D /usr/local/pgsql/data -l serverlog'
```

Here are a few more operating-system-specific suggestions. (In each case be sure to use the proper installation directory and user name where we show generic values.)

- For FreeBSD, look at the file [contrib/start-scripts/freebsd](#) in the PostgreSQL source distribution.
- On OpenBSD, add the following lines to the file `/etc/rc.local`:

```
if [ -x /usr/local/pgsql/bin/pg_ctl -a -x /usr/local/pgsql/bin/postgres ]; then
    su -l postgres -c '/usr/local/pgsql/bin/pg_ctl start -s -l /var/postgresql/log -
D /usr/local/pgsql/data'
    echo -n ' postgresql'
fi
```

- On Linux systems add

```
/usr/local/pgsql/bin/pg_ctl start -l logfile -D /usr/local/pgsql/data
to /etc/rc.d/rc.local or /etc/rc.local.
```

When using `systemd`, you can use the following service unit file (e.g., at `/etc/systemd/system/postgresql.service`):

```
[Unit]
Description=Postgres Pro database server
Documentation=man:postgres(1)
After=network-online.target
Wants=network-online.target

[Service]
Type=notify
User=postgres
ExecStart=/usr/local/pgsql/bin/postgres -D /usr/local/pgsql/data
ExecReload=/bin/kill -HUP $MAINPID
KillMode=mixed
KillSignal=SIGINT
TimeoutSec=infinity

[Install]
WantedBy=multi-user.target
```

Using `Type=notify` requires that the server binary was built with `configure --with-systemd`.

Consider carefully the timeout setting. `systemd` has a default timeout of 90 seconds as of this writing and will kill a process that does not report readiness within that time. But a Postgres Pro server that might have to perform crash recovery at startup could take much longer to become ready. The suggested value of `infinity` disables the timeout logic.

- On NetBSD, use either the FreeBSD or Linux start scripts, depending on preference.
- On Solaris, create a file called `/etc/init.d/postgresql` that contains the following line:

```
su - postgres -c "/usr/local/pgsql/bin/pg_ctl start -l logfile -D /usr/local/pgsql/data"
```

Then, create a symbolic link to it in `/etc/rc3.d` as `S99postgresql`.

While the server is running, its PID is stored in the file `postmaster.pid` in the data directory. This is used to prevent multiple server instances from running in the same data directory and can also be used for shutting down the server.

18.3.1. shardmand Examples

18.3.1.1. Configuring a shardmand Service

`shardmand` settings are usually specified in the `/etc/shardman/shardmand-cluster0.env` file. If you want `shardmand` to connect to an `etcd` cluster at hosts `n1-n3` using port 2379 and all `Shardman` services to use the `debug` log level, you can use the following `env` file:

```
SDM_STORE_ENDPOINTS=http://n1:2379,http://n2:2379,http://n3:2379
SDM_LOG_LEVEL=debug
```

Note that you need to restart `shardmand@cluster0` service to apply new settings from the `env` file. For more details see [Postgres Pro Shardman Start-up with shardmand](#).

18.3.1.2. Showing shardmand Logs

To look at `shardmand` logs, you can use a `journalctl` command:

```
$ journalctl -u shardmand@cluster0.service
```

18.3.1.3. Restarting Postgres Pro Shardman Services

You can restart all Postgres Pro `Shardman` services on a node using a `systemctl` command:

```
$ systemctl restart shardmand@cluster0.service
```

18.3.2. Server Start-up Failures

There are several common reasons the server might fail to start. Check the server's log file, or start it by hand (without redirecting standard output or standard error) and see what error messages appear. Below we explain some of the most common error messages in more detail.

```
LOG:   could not bind IPv4 address "127.0.0.1": Address already in use
HINT:   Is another postmaster already running on port 5432? If not, wait a few seconds
        and retry.
FATAL:  could not create any TCP/IP sockets
```

This usually means just what it suggests: you tried to start another server on the same port where one is already running. However, if the kernel error message is not `Address already in use` or some variant of that, there might be a different problem. For example, trying to start a server on a reserved port number might draw something like:

```
$ postgres -p 666
LOG:   could not bind IPv4 address "127.0.0.1": Permission denied
HINT:   Is another postmaster already running on port 666? If not, wait a few seconds
        and retry.
FATAL:  could not create any TCP/IP sockets
```

A message like:

```
FATAL:  could not create shared memory segment: Invalid argument
DETAIL: Failed system call was shmget(key=5440001, size=4011376640, 03600).
```

probably means your kernel's limit on the size of shared memory is smaller than the work area Postgres Pro is trying to create (4011376640 bytes in this example). This is only likely to happen if you have set `shared_memory_type` to `sysv`. In that case, you can try starting the server with a smaller-than-normal number of buffers ([shared_buffers](#)), or reconfigure your kernel to increase the allowed shared memory size. You might also see this message when trying to start multiple servers on the same machine, if their total space requested exceeds the kernel limit.

An error like:

```
FATAL:  could not create semaphores: No space left on device
DETAIL: Failed system call was semget(5440126, 17, 03600).
```

does *not* mean you've run out of disk space. It means your kernel's limit on the number of System V semaphores is smaller than the number Postgres Pro wants to create. As above, you might be able to work around the problem by starting the server with a reduced number of allowed connections ([max_connections](#)), but you'll eventually want to increase the kernel limit.

Details about configuring System V IPC facilities are given in [Section 18.4.1](#).

18.3.3. Client Connection Problems

Although the error conditions possible on the client side are quite varied and application-dependent, a few of them might be directly related to how the server was started. Conditions other than those shown below should be documented with the respective client application.

```
psql: error: connection to server at "server.joe.com" (123.123.123.123), port 5432
failed: Connection refused
        Is the server running on that host and accepting TCP/IP connections?
```

This is the generic “I couldn't find a server to talk to” failure. It looks like the above when TCP/IP communication is attempted. A common mistake is to forget to configure the server to allow TCP/IP connections.

Alternatively, you might get this when attempting Unix-domain socket communication to a local server:

```
psql: error: connection to server on socket "/tmp/.s.PGSQL.5432" failed: No such file
or directory
```

Is the server running locally and accepting connections on that socket?

If the server is indeed running, check that the client's idea of the socket path (here `/tmp`) agrees with the server's `unix_socket_directories` setting.

A connection failure message always shows the server address or socket path name, which is useful in verifying that the client is trying to connect to the right place. If there is in fact no server listening there, the kernel error message will typically be either `Connection refused` or `No such file or directory`, as illustrated. (It is important to realize that `Connection refused` in this context does *not* mean that the server got your connection request and rejected it. That case will produce a different message, as shown in [Section 20.15](#).) Other error messages such as `Connection timed out` might indicate more fundamental problems, like lack of network connectivity, or a firewall blocking the connection.

18.4. Managing Kernel Resources

Postgres Pro can sometimes exhaust various operating system resource limits, especially when multiple copies of the server are running on the same system, or in very large installations. This section explains the kernel resources used by Postgres Pro and the steps you can take to resolve problems related to kernel resource consumption.

18.4.1. Shared Memory and Semaphores

Postgres Pro requires the operating system to provide inter-process communication (IPC) features, specifically shared memory and semaphores. Unix-derived systems typically provide “System V” IPC, “POSIX” IPC, or both. Windows has its own implementation of these features and is not discussed here.

By default, Postgres Pro allocates a very small amount of System V shared memory, as well as a much larger amount of anonymous `mmap` shared memory. Alternatively, a single large System V shared memory region can be used (see [shared_memory_type](#)). In addition a significant number of semaphores, which can be either System V or POSIX style, are created at server startup. Currently, POSIX semaphores are used on Linux and FreeBSD systems while other platforms use System V semaphores.

System V IPC features are typically constrained by system-wide allocation limits. When Postgres Pro exceeds one of these limits, the server will refuse to start and should leave an instructive error message describing the problem and what to do about it. (See also [Section 18.3.2](#).) The relevant kernel parameters are named consistently across different systems; [Table 18.1](#) gives an overview. The methods to set them, however, vary. Suggestions for some platforms are given below.

Table 18.1. System V IPC Parameters

Name	Description	Values needed to run one Postgres Pro instance
SHMMAX	Maximum size of shared memory segment (bytes)	at least 1kB, but the default is usually much higher
SHMMIN	Minimum size of shared memory segment (bytes)	1
SHMALL	Total amount of shared memory available (bytes or pages)	same as SHMMAX if bytes, or <code>ceil(SHMMAX/PAGE_SIZE)</code> if pages, plus room for other applications
SHMSEG	Maximum number of shared memory segments per process	only 1 segment is needed, but the default is much higher
SHMMNI	Maximum number of shared memory segments system-wide	like SHMSEG plus room for other applications
SEMMNI	Maximum number of semaphore identifiers (i.e., sets)	at least <code>ceil((max_connections + autovacuum_max_workers + max_wal_senders + max_worker_processes + 7) / 19)</code> plus room for other applications

Name	Description	Values needed to run one Postgres Pro instance
SEMMNS	Maximum number of semaphores system-wide	$\text{ceil}((\text{max_connections} + \text{autovacuum_max_workers} + \text{max_wal_senders} + \text{max_worker_processes} + 7) / 19) * 20$ plus room for other applications
SEMMSL	Maximum number of semaphores per set	at least 20
SEMAP	Number of entries in semaphore map	see text
SEVMX	Maximum value of semaphore	at least 1000 (The default is often 32767; do not change unless necessary)

Postgres Pro requires a few bytes of System V shared memory (typically 48 bytes, on 64-bit platforms) for each copy of the server. On most modern operating systems, this amount can easily be allocated. However, if you are running many copies of the server or you explicitly configure the server to use large amounts of System V shared memory (see [shared_memory_type](#) and [dynamic_shared_memory_type](#)), it may be necessary to increase `SHMALL`, which is the total amount of System V shared memory system-wide. Note that `SHMALL` is measured in pages rather than bytes on many systems.

Less likely to cause problems is the minimum size for shared memory segments (`SHMMIN`), which should be at most approximately 32 bytes for Postgres Pro (it is usually just 1). The maximum number of segments system-wide (`SHMMNI`) or per-process (`SHMSEG`) are unlikely to cause a problem unless your system has them set to zero.

When using System V semaphores, Postgres Pro uses one semaphore per allowed connection ([max_connections](#)), allowed autovacuum worker process ([autovacuum_max_workers](#)), allowed WAL sender process ([max_wal_senders](#)), and allowed background process ([max_worker_processes](#)), in sets of 19. Each such set will also contain a 20th semaphore which contains a “magic number”, to detect collision with semaphore sets used by other applications. The maximum number of semaphores in the system is set by `SEMMNS`, which consequently must be at least as high as `max_connections` plus `autovacuum_max_workers` plus `max_wal_senders`, plus `max_worker_processes`, plus one extra for each 19 allowed connections plus workers (see the formula in [Table 18.1](#)). The parameter `SEMMNI` determines the limit on the number of semaphore sets that can exist on the system at one time. Hence this parameter must be at least $\text{ceil}((\text{max_connections} + \text{autovacuum_max_workers} + \text{max_wal_senders} + \text{max_worker_processes} + 7) / 19)$. Lowering the number of allowed connections is a temporary workaround for failures, which are usually confusingly worded “No space left on device”, from the function `semget`.

In some cases it might also be necessary to increase `SEMAP` to be at least on the order of `SEMMNS`. If the system has this parameter (many do not), it defines the size of the semaphore resource map, in which each contiguous block of available semaphores needs an entry. When a semaphore set is freed it is either added to an existing entry that is adjacent to the freed block or it is registered under a new map entry. If the map is full, the freed semaphores get lost (until reboot). Fragmentation of the semaphore space could over time lead to fewer available semaphores than there should be.

Various other settings related to “semaphore undo”, such as `SEMMNU` and `SEMUME`, do not affect Postgres Pro.

When using POSIX semaphores, the number of semaphores needed is the same as for System V, that is one semaphore per allowed connection ([max_connections](#)), allowed autovacuum worker process ([autovacuum_max_workers](#)), allowed WAL sender process ([max_wal_senders](#)), and allowed background process ([max_worker_processes](#)). On the platforms where this option is preferred, there is no specific kernel limit on the number of POSIX semaphores.

FreeBSD

The default shared memory settings are usually good enough, unless you have set `shared_memory_type` to `sysv`. System V semaphores are not used on this platform.

The default IPC settings can be changed using the `sysctl` or `loader` interfaces. The following parameters can be set using `sysctl`:

```
# sysctl kern.ipc.shmall=32768
# sysctl kern.ipc.shmmax=134217728
```

To make these settings persist over reboots, modify `/etc/sysctl.conf`.

If you have set `shared_memory_type` to `sysv`, you might also want to configure your kernel to lock System V shared memory into RAM and prevent it from being paged out to swap. This can be accomplished using the `sysctl` setting `kern.ipc.shm_use_phys`.

If running in a FreeBSD jail, you should set its `sysvshm` parameter to `new`, so that it has its own separate System V shared memory namespace. (Before FreeBSD 11.0, it was necessary to enable shared access to the host's IPC namespace from jails, and take measures to avoid collisions.)

NetBSD

The default shared memory settings are usually good enough, unless you have set `shared_memory_type` to `sysv`. You will usually want to increase `kern.ipc.semni` and `kern.ipc.semns`, as NetBSD's default settings for these are uncomfortably small.

IPC parameters can be adjusted using `sysctl`, for example:

```
# sysctl -w kern.ipc.semni=100
```

To make these settings persist over reboots, modify `/etc/sysctl.conf`.

If you have set `shared_memory_type` to `sysv`, you might also want to configure your kernel to lock System V shared memory into RAM and prevent it from being paged out to swap. This can be accomplished using the `sysctl` setting `kern.ipc.shm_use_phys`.

OpenBSD

The default shared memory settings are usually good enough, unless you have set `shared_memory_type` to `sysv`. You will usually want to increase `kern.semni` and `kern.semns`, as OpenBSD's default settings for these are uncomfortably small.

IPC parameters can be adjusted using `sysctl`, for example:

```
# sysctl kern.semni=100
```

To make these settings persist over reboots, modify `/etc/sysctl.conf`.

Linux

The default shared memory settings are usually good enough, unless you have set `shared_memory_type` to `sysv`, and even then only on older kernel versions that shipped with low defaults. System V semaphores are not used on this platform.

The shared memory size settings can be changed via the `sysctl` interface. For example, to allow 16 GB:

```
$ sysctl -w kernel.shmmax=17179869184
$ sysctl -w kernel.shmall=4194304
```

To make these settings persist over reboots, see `/etc/sysctl.conf`.

macOS

The default shared memory and semaphore settings are usually good enough, unless you have set `shared_memory_type` to `sysv`.

The recommended method for configuring shared memory in macOS is to create a file named `/etc/sysctl.conf`, containing variable assignments such as:

```
kern.sysv.shmmax=4194304
kern.sysv.shmmin=1
kern.sysv.shmmni=32
kern.sysv.shmseg=8
kern.sysv.shmall=1024
```

Note that in some macOS versions, *all five* shared-memory parameters must be set in `/etc/sysctl.conf`, else the values will be ignored.

SHMMAX can only be set to a multiple of 4096.

SHMALL is measured in 4 kB pages on this platform.

It is possible to change all but SHMMNI on the fly, using `sysctl`. But it's still best to set up your preferred values via `/etc/sysctl.conf`, so that the values will be kept across reboots.

Solaris illumos

The default shared memory and semaphore settings are usually good enough for most Postgres Pro applications. Solaris defaults to a SHMMAX of one-quarter of system RAM. To further adjust this setting, use a project setting associated with the `postgres` user. For example, run the following as `root`:

```
projadd -c "Postgres Pro DB User" -K "project.max-shm-memory=(privileged,8GB,deny) "
-U postgres -G postgres user.postgres
```

This command adds the `user.postgres` project and sets the shared memory maximum for the `postgres` user to 8GB, and takes effect the next time that user logs in, or when you restart Postgres Pro (not reload). The above assumes that Postgres Pro is run by the `postgres` user in the `postgres` group. No server reboot is required.

Other recommended kernel setting changes for database servers which will have a large number of connections are:

```
project.max-shm-ids=(priv,32768,deny)
project.max-sem-ids=(priv,4096,deny)
project.max-msg-ids=(priv,4096,deny)
```

Additionally, if you are running Postgres Pro inside a zone, you may need to raise the zone resource usage limits as well. See "Chapter2: Projects and Tasks" in the *System Administrator's Guide* for more information on `projects` and `prctl`.

18.4.2. systemd RemoveIPC

If `systemd` is in use, some care must be taken that IPC resources (including shared memory) are not prematurely removed by the operating system. This is especially of concern when installing Postgres Pro from source. Users of distribution packages of Postgres Pro are less likely to be affected, as the `postgres` user is then normally created as a system user.

The setting `RemoveIPC` in `logind.conf` controls whether IPC objects are removed when a user fully logs out. System users are exempt. This setting defaults to on in stock `systemd`, but some operating system distributions default it to off.

A typical observed effect when this setting is on is that shared memory objects used for parallel query execution are removed at apparently random times, leading to errors and warnings while attempting to open and remove them, like

```
WARNING: could not remove shared memory segment "/PostgreSQL.1450751626": No such file
or directory
```

Different types of IPC objects (shared memory vs. semaphores, System V vs. POSIX) are treated slightly differently by `systemd`, so one might observe that some IPC resources are not removed in the same way as others. But it is not advisable to rely on these subtle differences.

A “user logging out” might happen as part of a maintenance job or manually when an administrator logs in as the `postgres` user or something similar, so it is hard to prevent in general.

What is a “system user” is determined at `systemd` compile time from the `SYS_UID_MAX` setting in `/etc/login.defs`.

Packaging and deployment scripts should be careful to create the `postgres` user as a system user by using `useradd -r`, `adduser --system`, or equivalent.

Alternatively, if the user account was created incorrectly or cannot be changed, it is recommended to set

`RemoveIPC=no`

in `/etc/systemd/logind.conf` or another appropriate configuration file.

Caution

At least one of these two things has to be ensured, or the Postgres Pro server will be very unreliable.

18.4.3. Resource Limits

Unix-like operating systems enforce various kinds of resource limits that might interfere with the operation of your Postgres Pro server. Of particular importance are limits on the number of processes per user, the number of open files per process, and the amount of memory available to each process. Each of these have a “hard” and a “soft” limit. The soft limit is what actually counts but it can be changed by the user up to the hard limit. The hard limit can only be changed by the root user. The system call `setrlimit` is responsible for setting these parameters. The shell's built-in command `ulimit` (Bourne shells) or `limit` (csh) is used to control the resource limits from the command line. On BSD-derived systems the file `/etc/login.conf` controls the various resource limits set during login. See the operating system documentation for details. The relevant parameters are `maxproc`, `openfiles`, and `datasize`. For example:

```
default:\
...
      :datasize-cur=256M:\
      :maxproc-cur=256:\
      :openfiles-cur=256:\
...
```

(`-cur` is the soft limit. Append `-max` to set the hard limit.)

Kernels can also have system-wide limits on some resources.

- On Linux the kernel parameter `fs.file-max` determines the maximum number of open files that the kernel will support. It can be changed with `sysctl -w fs.file-max=N`. To make the setting persist across reboots, add an assignment in `/etc/sysctl.conf`. The maximum limit of files per process is fixed at the time the kernel is compiled; see `/usr/src/linux/Documentation/proc.txt` for more information.

The Postgres Pro server uses one process per connection so you should provide for at least as many processes as allowed connections, in addition to what you need for the rest of your system. This is usually not a problem but if you run several servers on one machine things might get tight.

The factory default limit on open files is often set to “socially friendly” values that allow many users to coexist on a machine without using an inappropriate fraction of the system resources. If you run many servers on a machine this is perhaps what you want, but on dedicated servers you might want to raise this limit.

On the other side of the coin, some systems allow individual processes to open large numbers of files; if more than a few processes do so then the system-wide limit can easily be exceeded. If you

find this happening, and you do not want to alter the system-wide limit, you can set Postgres Pro's `max_files_per_process` configuration parameter to limit the consumption of open files.

Another kernel limit that may be of concern when supporting large numbers of client connections is the maximum socket connection queue length. If more than that many connection requests arrive within a very short period, some may get rejected before the Postgres Pro server can service the requests, with those clients receiving unhelpful connection failure errors such as “Resource temporarily unavailable” or “Connection refused”. The default queue length limit is 128 on many platforms. To raise it, adjust the appropriate kernel parameter via `sysctl`, then restart the Postgres Pro server. The parameter is variously named `net.core.somaxconn` on Linux, `kern.ipc.soacceptqueue` on newer FreeBSD, and `kern.ipc.somaxconn` on macOS and other BSD variants.

18.4.4. Linux Memory Overcommit

The default virtual memory behavior on Linux is not optimal for Postgres Pro. Because of the way that the kernel implements memory overcommit, the kernel might terminate the Postgres Pro postmaster (the supervisor server process) if the memory demands of either Postgres Pro or another process cause the system to run out of virtual memory.

If this happens, you will see a kernel message that looks like this (consult your system documentation and configuration on where to look for such a message):

```
Out of Memory: Killed process 12345 (postgres).
```

This indicates that the `postgres` process has been terminated due to memory pressure. Although existing database connections will continue to function normally, no new connections will be accepted. To recover, Postgres Pro will need to be restarted.

One way to avoid this problem is to run Postgres Pro on a machine where you can be sure that other processes will not run the machine out of memory. If memory is tight, increasing the swap space of the operating system can help avoid the problem, because the out-of-memory (OOM) killer is invoked only when physical memory and swap space are exhausted.

If Postgres Pro itself is the cause of the system running out of memory, you can avoid the problem by changing your configuration. In some cases, it may help to lower memory-related configuration parameters, particularly `shared_buffers`, `work_mem`, and `hash_mem_multiplier`. In other cases, the problem may be caused by allowing too many connections to the database server itself. In many cases, it may be better to reduce `max_connections` and instead make use of external connection-pooling software.

It is possible to modify the kernel's behavior so that it will not “overcommit” memory. Although this setting will not prevent the *OOM killer* from being invoked altogether, it will lower the chances significantly and will therefore lead to more robust system behavior. This is done by selecting strict overcommit mode via `sysctl`:

```
sysctl -w vm.overcommit_memory=2
```

or placing an equivalent entry in `/etc/sysctl.conf`. You might also wish to modify the related setting `vm.overcommit_ratio`. For details see the kernel documentation file <https://www.kernel.org/doc/Documentation/vm/overcommit-accounting>.

Another approach, which can be used with or without altering `vm.overcommit_memory`, is to set the process-specific *OOM score adjustment* value for the postmaster process to `-1000`, thereby guaranteeing it will not be targeted by the OOM killer. The simplest way to do this is to execute

```
echo -1000 > /proc/self/oom_score_adj
```

in the Postgres Pro startup script just before invoking `postgres`. Note that this action must be done as root, or it will have no effect; so a root-owned startup script is the easiest place to do it. If you do this, you should also set these environment variables in the startup script before invoking `postgres`:

```
export PG_OOM_ADJUST_FILE=/proc/self/oom_score_adj
export PG_OOM_ADJUST_VALUE=0
```

These settings will cause postmaster child processes to run with the normal OOM score adjustment of zero, so that the OOM killer can still target them at need. You could use some other value for `PG_OOM_ADJUST_VALUE` if you want the child processes to run with some other OOM score adjustment. (`PG_OOM_ADJUST_VALUE` can also be omitted, in which case it defaults to zero.) If you do not set `PG_OOM_ADJUST_FILE`, the child processes will run with the same OOM score adjustment as the postmaster, which is unwise since the whole point is to ensure that the postmaster has a preferential setting.

18.4.5. Linux Huge Pages

Using huge pages reduces overhead when using large contiguous chunks of memory, as Postgres Pro does, particularly when using large values of [shared_buffers](#). To use this feature in Postgres Pro you need a kernel with `CONFIG_HUGETLBFS=y` and `CONFIG_HUGETLB_PAGE=y`. You will also have to configure the operating system to provide enough huge pages of the desired size. The runtime-computed parameter [shared_memory_size_in_huge_pages](#) reports the number of huge pages required. This parameter can be viewed before starting the server with a `postgres` command like:

```
$ postgres -D $PGDATA -C shared_memory_size_in_huge_pages
3170
$ grep ^Hugepagesize /proc/meminfo
Hugepagesize:          2048 kB
$ ls /sys/kernel/mm/hugepages
hugepages-1048576kB  hugepages-2048kB
```

In this example the default is 2MB, but you can also explicitly request either 2MB or 1GB with [huge_page_size](#) to adapt the number of pages calculated by `shared_memory_size_in_huge_pages`. While we need at least 3170 huge pages in this example, a larger setting would be appropriate if other programs on the machine also need huge pages. We can set this with:

```
# sysctl -w vm.nr_hugepages=3170
```

Don't forget to add this setting to `/etc/sysctl.conf` so that it is reapplied after reboots. For non-default huge page sizes, we can instead use:

```
# echo 3170 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
```

It is also possible to provide these settings at boot time using kernel parameters such as `hugepagesz=2M` `hugepages=3170`.

Sometimes the kernel is not able to allocate the desired number of huge pages immediately due to fragmentation, so it might be necessary to repeat the command or to reboot. (Immediately after a reboot, most of the machine's memory should be available to convert into huge pages.) To verify the huge page allocation situation for a given size, use:

```
$ cat /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
```

It may also be necessary to give the database server's operating system user permission to use huge pages by setting `vm.hugetlb_shm_group` via `sysctl`, and/or give permission to lock memory with `ulimit -l`.

The default behavior for huge pages in Postgres Pro is to use them when possible, with the system's default huge page size, and to fall back to normal pages on failure. To enforce the use of huge pages, you can set [huge_pages](#) to `on` in `postgresql.conf`. Note that with this setting Postgres Pro will fail to start if not enough huge pages are available.

For a detailed description of the Linux huge pages feature have a look at <https://www.kernel.org/doc/Documentation/vm/hugetlbpage.txt>.

18.4.6. Resource Prioritization

Postgres Pro Shardman provides an experimental feature for resource prioritization.

On systems with limited resources or under heavy load, you may need to prioritize transaction execution, so that some transactions are executed more quickly than the other. For example, you may want to

execute simple user queries as fast as possible, even if it delays less urgent tasks, such as complex OLAP queries that may be running at the same time. Postgres Pro Shardman enables you to configure resource prioritization policy, which can slow down a particular session based on the amount of CPU, I/O read, and I/O write resources this session consumes as compared to other sessions.

By default, resource prioritization is disabled, so all backends have equal access to all the available resources. You can assign weight to each backend to control the amount of resources each session can use within the specified time interval. Depending on the current resource consumption, Postgres Pro Shardman may suspend backends with lower weight from time to time to ensure that high-priority sessions have more resources available.

To enable prioritization in your database cluster:

1. Configure the time interval for collecting usage statistics for all active backends by setting the `usage_tracking_interval` parameter in the `postgresql.conf` file and reload the server configuration.

Once the `usage_tracking_interval` parameter is set, Postgres Pro Shardman starts collecting statistics on resource usage at the specified interval.

Tip

Avoid setting `usage_tracking_interval` to small values as frequent statistics collection can cause overhead.

2. Depending on the resources you need to control, modify one or more of the following parameters for the sessions you would like to prioritize:
 - `session_cpu_weight` — CPU usage.
 - `session_ioread_weight` — I/O read throughput.
 - `session_iowrite_weight` — I/O write throughput.

These parameters can take weight values 1, 2, 4, and 8. The higher the value, the more resources the session can use. Sessions with the same weight have the same priority for resource usage, so if you assign equal weights to all sessions, performance is not affected, regardless of the weight value. By default, all sessions have weight 4 for all types of resources.

For all possible ways of modifying configuration for a particular session, see [Section 19.1](#).

Once you change the weight of one or more sessions, Postgres Pro Shardman enables prioritization policy based on the assigned weight values and the usage statistics measured for the previous `usage_tracking_interval`. Thus, session activity is adjusted for each usage tracking interval, if required.

Note

Even though weights for each resource are assigned separately, prioritizing a session by one resource can indirectly affect the session performance with regard to other resources.

18.5. Shutting Down the Server

For distributed system specifics, refer to the `shardmanctl` tool.

There are several ways to shut down the database server. Under the hood, they all reduce to sending a signal to the supervisor `postgres` process.

If you are using a pre-packaged version of PostgreSQL, and you used its provisions for starting the server, then you should also use its provisions for stopping the server. Consult the package-level documentation for details.

When managing the server directly, you can control the type of shutdown by sending different signals to the `postgres` process:

SIGTERM

This is the *Smart Shutdown* mode. After receiving SIGTERM, the server disallows new connections, but lets existing sessions end their work normally. It shuts down only after all of the sessions terminate. If the server is in recovery when a smart shutdown is requested, recovery and streaming replication will be stopped only after all regular sessions have terminated.

SIGINT

This is the *Fast Shutdown* mode. The server disallows new connections and sends all existing server processes SIGTERM, which will cause them to abort their current transactions and exit promptly. It then waits for all server processes to exit and finally shuts down.

SIGQUIT

This is the *Immediate Shutdown* mode. The server will send SIGQUIT to all child processes and wait for them to terminate. If any do not terminate within 5 seconds, they will be sent SIGKILL. The supervisor server process exits as soon as all child processes have exited, without doing normal database shutdown processing. This will lead to recovery (by replaying the WAL log) upon next start-up. This is recommended only in emergencies.

The `pg_ctl` program provides a convenient interface for sending these signals to shut down the server. Alternatively, you can send the signal directly using `kill` on non-Windows systems. The PID of the `postgres` process can be found using the `ps` program, or from the file `postmaster.pid` in the data directory. For example, to do a fast shutdown:

```
$ kill -INT `head -1 /usr/local/pgsql/data/postmaster.pid`
```

Important

It is best not to use SIGKILL to shut down the server. Doing so will prevent the server from releasing shared memory and semaphores. Furthermore, SIGKILL kills the `postgres` process without letting it relay the signal to its subprocesses, so it might be necessary to kill the individual subprocesses by hand as well.

To terminate an individual session while allowing other sessions to continue, use `pg_terminate_backend()` (see [Table 9.94](#)) or send a SIGTERM signal to the child process associated with the session.

18.6. Upgrading a Postgres Pro Cluster

This section discusses how to upgrade your database from one Postgres Pro Shardman release to a newer one. It is best to review the [Release Notes](#) before an upgrade and look for any changes that may cause issues for your application. You can proceed to upgrade if there are no potential issues.

The process of updating a Postgres Pro Shardman consists of several steps that must be performed sequentially:

1. Upgrade Postgres Pro Shardman packages.
2. Restart all Postgres Pro Shardman services and database instances.
3. Upgrade database shardman extension.

18.6.1. Upgrade Packages

18.6.1.1. APT-based Systems

To upgrade packages, typically run the following command:

```
$ apt update && apt --only-upgrade install shardman-tools shardman-services  
postgrespro-sdm-17-contrib postgrespro-sdm-17-server
```

or upgrade all packages:

```
$ apt update && apt upgrade
```

Check that all packages have been updated on each node:

```
$ dpkg -l | grep -E ' (postgrespro|shardman) '
```

18.6.1.2. RPM-based systems

To upgrade packages, typically run the following command:

```
$ yum update shardman-tools shardman-services postgrespro-sdm-17-contrib postgrespro-  
sdm-17-server
```

or upgrade all packages:

```
$ yum update
```

Check that all packages have been updated on each node:

```
$ yum list --installed | grep -E ' (postgrespro|shardman) '
```

18.6.2. Restart Postgres Pro Shardman Services and Database Instances

After updating the packages, you need to restart all cluster services. It can be done with a single `shardmanctl restart` command:

```
$ shardmanctl --cluster-name cluster0 --store-endpoints http://etcd1:2379,http://  
etcd2:2379,http://etcd3:2379 restart
```

You can skip the `--cluster-name` and `--store-endpoints` options by setting the `SDM_CLUSTER_NAME` and `SDM_STORE_ENDPOINTS` environment variables as in the example below:

```
export SDM_STORE_ENDPOINTS=http://etcd1:2379,http://  
etcd2:2379,http://etcd3:2379  
export SDM_CLUSTER_NAME=cluster0
```

18.6.3. Upgrade the Extension

After restarting services of the cluster, you should update the server extensions by running the following command:

```
$ shardmanctl --cluster-name cluster0 --store-endpoints http://etcd1:2379,http://  
etcd2:2379,http://etcd3:2379 upgrade
```

In the case when the shardman extension version and server library version are different, distributed queries and Postgres Pro Shardman DDL will not work.

Postgres Pro Shardman extensions try to ensure that they do not communicate with incompatible software. Incompatibilities can arise for several reasons: the shardman shared library version does not match the extension version or the remote server version does not match the local server version. In case when the extension and library versions mismatch, Postgres Pro Shardman cannot modify its metadata and will refuse to perform operations on global objects until the extension is updated. In case when the remote server version does not match the local server version or when they belong to different clusters, Postgres Pro Shardman will refuse to communicate with the server.

18.6.4. Upgrading Data via Replication

It is also possible to use logical replication methods to create a standby server with the updated version of Postgres Pro. This is possible because logical replication supports replication between different major

versions of Postgres Pro. The standby can be on the same computer or a different computer. Once it has synced up with the primary server (running the older version of Postgres Pro), you can switch primaries and make the standby the primary and shut down the older database instance. Such a switch-over results in only several seconds of downtime for an upgrade.

This method of upgrading can be performed using the built-in logical replication facilities as well as using external logical replication systems such as pglogical, Slony, Londiste, and Bucardo.

18.7. Preventing Server Spoofing

While the server is running, it is not possible for a malicious user to take the place of the normal database server. However, when the server is down, it is possible for a local user to spoof the normal server by starting their own server. The spoof server could read passwords and queries sent by clients, but could not return any data because the `PGDATA` directory would still be secure because of directory permissions. Spoofing is possible because any user can start a database server; a client cannot identify an invalid server unless it is specially configured.

One way to prevent spoofing of `local` connections is to use a Unix domain socket directory ([unix_socket_directories](#)) that has write permission only for a trusted local user. This prevents a malicious user from creating their own socket file in that directory. If you are concerned that some applications might still reference `/tmp` for the socket file and hence be vulnerable to spoofing, during operating system startup create a symbolic link `/tmp/.s.PGSQL.5432` that points to the relocated socket file. You also might need to modify your `/tmp` cleanup script to prevent removal of the symbolic link.

Another option for `local` connections is for clients to use [requirepeer](#) to specify the required owner of the server process connected to the socket.

To prevent spoofing on TCP connections, either use SSL certificates and make sure that clients check the server's certificate, or use GSSAPI encryption (or both, if they're on separate connections).

To prevent spoofing with SSL, the server must be configured to accept only `hostssl` connections ([Section 20.1](#)) and have SSL key and certificate files ([Section 18.9](#)). The TCP client must connect using `sslmode=verify-ca` or `verify-full` and have the appropriate root certificate file installed ([Section 35.19.1](#)). Alternatively the [system CA pool](#), as defined by the SSL implementation, can be used using `sslrootcert=system`; in this case, `sslmode=verify-full` is forced for safety, since it is generally trivial to obtain certificates which are signed by a public CA.

To prevent server spoofing from occurring when using [scram-sha-256](#) password authentication over a network, you should ensure that you connect to the server using SSL and with one of the anti-spoofing methods described in the previous paragraph. Additionally, the SCRAM implementation in libpq cannot protect the entire authentication exchange, but using the `channel_binding=require` connection parameter provides a mitigation against server spoofing. An attacker that uses a rogue server to intercept a SCRAM exchange can use offline analysis to potentially determine the hashed password from the client.

To prevent spoofing with GSSAPI, the server must be configured to accept only `hostgssenc` connections ([Section 20.1](#)) and use `gss` authentication with them. The TCP client must connect using `gssencmode=require`.

18.8. Encryption Options

Postgres Pro offers encryption at several levels, and provides flexibility in protecting data from disclosure due to database server theft, unscrupulous administrators, and insecure networks. Encryption might also be required to secure sensitive data such as medical records or financial transactions.

Password Encryption

Database user passwords are stored as hashes (determined by the setting [password_encryption](#)), so the administrator cannot determine the actual password assigned to the user. If SCRAM or MD5 encryption is used for client authentication, the unencrypted password is never even temporarily present on the server because the client encrypts it before being sent across the network. SCRAM

is preferred, because it is an Internet standard and is more secure than the Postgres Pro-specific MD5 authentication protocol.

Encryption For Specific Columns

The `pgcrypto` module allows certain fields to be stored encrypted. This is useful if only some of the data is sensitive. The client supplies the decryption key and the data is decrypted on the server and then sent to the client.

The decrypted data and the decryption key are present on the server for a brief time while it is being decrypted and communicated between the client and server. This presents a brief moment where the data and keys can be intercepted by someone with complete access to the database server, such as the system administrator.

Data Partition Encryption

Storage encryption can be performed at the file system level or the block level. Linux file system encryption options include eCryptfs and EncFS, while FreeBSD uses PEFS. Block level or full disk encryption options include dm-crypt + LUKS on Linux and GEOM modules geli and gbde on FreeBSD. Many other operating systems support this functionality, including Windows.

This mechanism prevents unencrypted data from being read from the drives if the drives or the entire computer is stolen. This does not protect against attacks while the file system is mounted, because when mounted, the operating system provides an unencrypted view of the data. However, to mount the file system, you need some way for the encryption key to be passed to the operating system, and sometimes the key is stored somewhere on the host that mounts the disk.

Encrypting Data Across A Network

SSL connections encrypt all data sent across the network: the password, the queries, and the data returned. The `pg_hba.conf` file allows administrators to specify which hosts can use non-encrypted connections (`host`) and which require SSL-encrypted connections (`hostssl`). Also, clients can specify that they connect to servers only via SSL.

GSSAPI-encrypted connections encrypt all data sent across the network, including queries and data returned. (No password is sent across the network.) The `pg_hba.conf` file allows administrators to specify which hosts can use non-encrypted connections (`host`) and which require GSSAPI-encrypted connections (`hostgssenc`). Also, clients can specify that they connect to servers only on GSSAPI-encrypted connections (`gssencmode=require`).

Stunnel or SSH can also be used to encrypt transmissions.

SSL Host Authentication

It is possible for both the client and server to provide SSL certificates to each other. It takes some extra configuration on each side, but this provides stronger verification of identity than the mere use of passwords. It prevents a computer from pretending to be the server just long enough to read the password sent by the client. It also helps prevent “man in the middle” attacks where a computer between the client and server pretends to be the server and reads and passes all data between the client and server.

Client-Side Encryption

If the system administrator for the server's machine cannot be trusted, it is necessary for the client to encrypt the data; this way, unencrypted data never appears on the database server. Data is encrypted on the client before being sent to the server, and database results have to be decrypted on the client before being used.

18.9. Secure TCP/IP Connections with SSL

Postgres Pro has native support for using SSL connections to encrypt client/server communications for increased security.

The terms SSL and TLS are often used interchangeably to mean a secure encrypted connection using a TLS protocol. SSL protocols are the precursors to TLS protocols, and the term SSL is still used for encrypted connections even though SSL protocols are no longer supported. SSL is used interchangeably with TLS in Postgres Pro.

18.9.1. Basic Setup

With SSL support compiled in, the Postgres Pro server can be started with support for encrypted connections using TLS protocols enabled by setting the parameter `ssl` to `on` in `postgresql.conf`. The server will listen for both normal and SSL connections on the same TCP port, and will negotiate with any connecting client on whether to use SSL. By default, this is at the client's option; see [Section 20.1](#) about how to set up the server to require use of SSL for some or all connections.

To start in SSL mode, files containing the server certificate and private key must exist. By default, these files are expected to be named `server.crt` and `server.key`, respectively, in the server's data directory, but other names and locations can be specified using the configuration parameters `ssl_cert_file` and `ssl_key_file`.

On Unix systems, the permissions on `server.key` must disallow any access to world or group; achieve this by the command `chmod 0600 server.key`. Alternatively, the file can be owned by root and have group read access (that is, 0640 permissions). That setup is intended for installations where certificate and key files are managed by the operating system. The user under which the Postgres Pro server runs should then be made a member of the group that has access to those certificate and key files.

If the data directory allows group read access then certificate files may need to be located outside of the data directory in order to conform to the security requirements outlined above. Generally, group access is enabled to allow an unprivileged user to backup the database, and in that case the backup software will not be able to read the certificate files and will likely error.

If the private key is protected with a passphrase, the server will prompt for the passphrase and will not start until it has been entered. Using a passphrase by default disables the ability to change the server's SSL configuration without a server restart, but see [ssl_passphrase_command_supports_reload](#). Furthermore, passphrase-protected private keys cannot be used at all on Windows.

The first certificate in `server.crt` must be the server's certificate because it must match the server's private key. The certificates of “intermediate” certificate authorities can also be appended to the file. Doing this avoids the necessity of storing intermediate certificates on clients, assuming the root and intermediate certificates were created with `v3_ca` extensions. (This sets the certificate's basic constraint of CA to `true`.) This allows easier expiration of intermediate certificates.

It is not necessary to add the root certificate to `server.crt`. Instead, clients must have the root certificate of the server's certificate chain.

18.9.2. OpenSSL Configuration

Postgres Pro reads the system-wide OpenSSL configuration file. By default, this file is named `openssl.cnf` and is located in the directory reported by `openssl version -d`. This default can be overridden by setting environment variable `OPENSSL_CONF` to the name of the desired configuration file.

OpenSSL supports a wide range of ciphers and authentication algorithms, of varying strength. While a list of ciphers can be specified in the OpenSSL configuration file, you can specify ciphers specifically for use by the database server by modifying `ssl_ciphers` in `postgresql.conf`.

Note

It is possible to have authentication without encryption overhead by using `NULL-SHA` or `NULL-MD5` ciphers. However, a man-in-the-middle could read and pass communications between client and server. Also, encryption overhead is minimal compared to the overhead of authentication. For these reasons NULL ciphers are not recommended.

18.9.3. Using Client Certificates

To require the client to supply a trusted certificate, place certificates of the root certificate authorities (CAs) you trust in a file in the data directory, set the parameter `ssl_ca_file` in `postgresql.conf` to the new file name, and add the authentication option `clientcert=verify-ca` or `clientcert=verify-full` to the appropriate `hostssl` line(s) in `pg_hba.conf`. A certificate will then be requested from the client during SSL connection startup. (See [Section 35.19](#) for a description of how to set up certificates on the client.)

For a `hostssl` entry with `clientcert=verify-ca`, the server will verify that the client's certificate is signed by one of the trusted certificate authorities. If `clientcert=verify-full` is specified, the server will not only verify the certificate chain, but it will also check whether the username or its mapping matches the `cn` (Common Name) of the provided certificate. Note that certificate chain validation is always ensured when the `cert` authentication method is used (see [Section 20.12](#)).

Intermediate certificates that chain up to existing root certificates can also appear in the `ssl_ca_file` file if you wish to avoid storing them on clients (assuming the root and intermediate certificates were created with `v3_ca` extensions). Certificate Revocation List (CRL) entries are also checked if the parameter `ssl_crl_file` or `ssl_crl_dir` is set.

The `clientcert` authentication option is available for all authentication methods, but only in `pg_hba.conf` lines specified as `hostssl`. When `clientcert` is not specified, the server verifies the client certificate against its CA file only if a client certificate is presented and the CA is configured.

There are two approaches to enforce that users provide a certificate during login.

The first approach makes use of the `cert` authentication method for `hostssl` entries in `pg_hba.conf`, such that the certificate itself is used for authentication while also providing ssl connection security. See [Section 20.12](#) for details. (It is not necessary to specify any `clientcert` options explicitly when using the `cert` authentication method.) In this case, the `cn` (Common Name) provided in the certificate is checked against the user name or an applicable mapping.

The second approach combines any authentication method for `hostssl` entries with the verification of client certificates by setting the `clientcert` authentication option to `verify-ca` or `verify-full`. The former option only enforces that the certificate is valid, while the latter also ensures that the `cn` (Common Name) in the certificate matches the user name or an applicable mapping.

18.9.4. SSL Server File Usage

[Table 18.2](#) summarizes the files that are relevant to the SSL setup on the server. (The shown file names are default names. The locally configured names could be different.)

Table 18.2. SSL Server File Usage

File	Contents	Effect
<code>ssl_cert_file</code> (\$PGDATA/server.crt)	server certificate	sent to client to indicate server's identity
<code>ssl_key_file</code> (\$PGDATA/server.key)	server private key	proves server certificate was sent by the owner; does not indicate certificate owner is trustworthy
<code>ssl_ca_file</code>	trusted certificate authorities	checks that client certificate is signed by a trusted certificate authority
<code>ssl_crl_file</code>	certificates revoked by certificate authorities	client certificate must not be on this list

The server reads these files at server start and whenever the server configuration is reloaded. On Windows systems, they are also re-read whenever a new backend process is spawned for a new client connection.

If an error in these files is detected at server start, the server will refuse to start. But if an error is detected during a configuration reload, the files are ignored and the old SSL configuration continues to be used. On Windows systems, if an error in these files is detected at backend start, that backend will be unable to establish an SSL connection. In all these cases, the error condition is reported in the server log.

18.9.5. Creating Certificates

To create a simple self-signed certificate for the server, valid for 365 days, use the following OpenSSL command, replacing *dbhost.yourdomain.com* with the server's host name:

```
openssl req -new -x509 -days 365 -nodes -text -out server.crt \  
-keyout server.key -subj "/CN=dbhost.yourdomain.com"
```

Then do:

```
chmod og-rwx server.key
```

because the server will reject the file if its permissions are more liberal than this. For more details on how to create your server private key and certificate, refer to the OpenSSL documentation.

While a self-signed certificate can be used for testing, a certificate signed by a certificate authority (CA) (usually an enterprise-wide root CA) should be used in production.

To create a server certificate whose identity can be validated by clients, first create a certificate signing request (CSR) and a public/private key file:

```
openssl req -new -nodes -text -out root.csr \  
-keyout root.key -subj "/CN=root.yourdomain.com"  
chmod og-rwx root.key
```

Then, sign the request with the key to create a root certificate authority (using the default OpenSSL configuration file location on Linux):

```
openssl x509 -req -in root.csr -text -days 3650 \  
-extfile /etc/ssl/openssl.cnf -extensions v3_ca \  
-signkey root.key -out root.crt
```

Finally, create a server certificate signed by the new root certificate authority:

```
openssl req -new -nodes -text -out server.csr \  
-keyout server.key -subj "/CN=dbhost.yourdomain.com"  
chmod og-rwx server.key
```

```
openssl x509 -req -in server.csr -text -days 365 \  
-CA root.crt -CAkey root.key -CAcreateserial \  
-out server.crt
```

server.crt and *server.key* should be stored on the server, and *root.crt* should be stored on the client so the client can verify that the server's leaf certificate was signed by its trusted root certificate. *root.key* should be stored offline for use in creating future certificates.

It is also possible to create a chain of trust that includes intermediate certificates:

```
# root  
openssl req -new -nodes -text -out root.csr \  
-keyout root.key -subj "/CN=root.yourdomain.com"  
chmod og-rwx root.key  
openssl x509 -req -in root.csr -text -days 3650 \  
-extfile /etc/ssl/openssl.cnf -extensions v3_ca \  
-signkey root.key -out root.crt  
  
# intermediate  
openssl req -new -nodes -text -out intermediate.csr \  
-keyout intermediate.key -subj "/CN=intermediate.yourdomain.com"
```

```
chmod og-rwx intermediate.key
openssl x509 -req -in intermediate.csr -text -days 1825 \
  -extfile /etc/ssl/openssl.cnf -extensions v3_ca \
  -CA root.crt -CAkey root.key -CAcreateserial \
  -out intermediate.crt

# leaf
openssl req -new -nodes -text -out server.csr \
  -keyout server.key -subj "/CN=dbhost.yourdomain.com"
chmod og-rwx server.key
openssl x509 -req -in server.csr -text -days 365 \
  -CA intermediate.crt -CAkey intermediate.key -CAcreateserial \
  -out server.crt
```

`server.crt` and `intermediate.crt` should be concatenated into a certificate file bundle and stored on the server. `server.key` should also be stored on the server. `root.crt` should be stored on the client so the client can verify that the server's leaf certificate was signed by a chain of certificates linked to its trusted root certificate. `root.key` and `intermediate.key` should be stored offline for use in creating future certificates.

18.10. Secure TCP/IP Connections with GSSAPI Encryption

Postgres Pro also has native support for using GSSAPI to encrypt client/server communications for increased security. Support requires that a GSSAPI implementation (such as MIT Kerberos) is installed on both client and server systems, and that support in Postgres Pro is enabled at build time.

18.10.1. Basic Setup

The Postgres Pro server will listen for both normal and GSSAPI-encrypted connections on the same TCP port, and will negotiate with any connecting client whether to use GSSAPI for encryption (and for authentication). By default, this decision is up to the client (which means it can be downgraded by an attacker); see [Section 20.1](#) about setting up the server to require the use of GSSAPI for some or all connections.

When using GSSAPI for encryption, it is common to use GSSAPI for authentication as well, since the underlying mechanism will determine both client and server identities (according to the GSSAPI implementation) in any case. But this is not required; another Postgres Pro authentication method can be chosen to perform additional verification.

Other than configuration of the negotiation behavior, GSSAPI encryption requires no setup beyond that which is necessary for GSSAPI authentication. (For more information on configuring that, see [Section 20.6](#).)

18.11. Secure TCP/IP Connections with SSH Tunnels

It is possible to use SSH to encrypt the network connection between clients and a Postgres Pro server. Done properly, this provides an adequately secure network connection, even for non-SSL-capable clients.

First make sure that an SSH server is running properly on the same machine as the Postgres Pro server and that you can log in using `ssh` as some user; you then can establish a secure tunnel to the remote server. A secure tunnel listens on a local port and forwards all traffic to a port on the remote machine. Traffic sent to the remote port can arrive on its `localhost` address, or different bind address if desired; it does not appear as coming from your local machine. This command creates a secure tunnel from the client machine to the remote machine `foo.com`:

```
ssh -L 63333:localhost:5432 joe@foo.com
```

The first number in the `-L` argument, 63333, is the local port number of the tunnel; it can be any unused port. (IANA reserves ports 49152 through 65535 for private use.) The name or IP address after this is the

remote bind address you are connecting to, i.e., `localhost`, which is the default. The second number, 5432, is the remote end of the tunnel, e.g., the port number your database server is using. In order to connect to the database server using this tunnel, you connect to port 63333 on the local machine:

```
psql -h localhost -p 63333 postgres
```

To the database server it will then look as though you are user `joe` on host `foo.com` connecting to the `localhost` bind address, and it will use whatever authentication procedure was configured for connections by that user to that bind address. Note that the server will not think the connection is SSL-encrypted, since in fact it is not encrypted between the SSH server and the Postgres Pro server. This should not pose any extra security risk because they are on the same machine.

In order for the tunnel setup to succeed you must be allowed to connect via `ssh` as `joe@foo.com`, just as if you had attempted to use `ssh` to create a terminal session.

You could also have set up port forwarding as

```
ssh -L 63333:foo.com:5432 joe@foo.com
```

but then the database server will see the connection as coming in on its `foo.com` bind address, which is not opened by the default setting `listen_addresses = 'localhost'`. This is usually not what you want.

If you have to “hop” to the database server via some login host, one possible setup could look like this:

```
ssh -L 63333:db.foo.com:5432 joe@shell.foo.com
```

Note that this way the connection from `shell.foo.com` to `db.foo.com` will not be encrypted by the SSH tunnel. SSH offers quite a few configuration possibilities when the network is restricted in various ways. Please refer to the SSH documentation for details.

Tip

Several other applications exist that can provide secure tunnels using a procedure similar in concept to the one just described.

18.12. Registering Event Log on Windows

To register a Windows event log library with the operating system, issue this command:

```
regsvr32 pgsql_library_directory/pgevent.dll
```

This creates registry entries used by the event viewer, under the default event source named `Postgres Pro`.

To specify a different event source name (see [event_source](#)), use the `/n` and `/i` options:

```
regsvr32 /n /i:event_source_name pgsql_library_directory/pgevent.dll
```

To unregister the event log library from the operating system, issue this command:

```
regsvr32 /u [/i:event_source_name] pgsql_library_directory/pgevent.dll
```

Note

To enable event logging in the database server, modify [log_destination](#) to include `eventlog` in `postgresql.conf`.

Chapter 19. Server Configuration

There are many configuration parameters that affect the behavior of the database system. In the first section of this chapter we describe how to interact with configuration parameters. The subsequent sections discuss each parameter in detail.

19.1. Setting Parameters

19.1.1. Parameter Names and Values

All parameter names are case-insensitive. Every parameter takes a value of one of five types: boolean, string, integer, floating point, or enumerated (enum). The type determines the syntax for setting the parameter:

- *Boolean*: Values can be written as `on`, `off`, `true`, `false`, `yes`, `no`, `1`, `0` (all case-insensitive) or any unambiguous prefix of one of these.
- *String*: In general, enclose the value in single quotes, doubling any single quotes within the value. Quotes can usually be omitted if the value is a simple number or identifier, however. (Values that match an SQL keyword require quoting in some contexts.)
- *Numeric (integer and floating point)*: Numeric parameters can be specified in the customary integer and floating-point formats; fractional values are rounded to the nearest integer if the parameter is of integer type. Integer parameters additionally accept hexadecimal input (beginning with `0x`) and octal input (beginning with `0`), but these formats cannot have a fraction. Do not use thousands separators. Quotes are not required, except for hexadecimal input.
- *Numeric with Unit*: Some numeric parameters have an implicit unit, because they describe quantities of memory or time. The unit might be bytes, kilobytes, blocks (typically eight kilobytes), milliseconds, seconds, or minutes. An unadorned numeric value for one of these settings will use the setting's default unit, which can be learned from `pg_settings.unit`. For convenience, settings can be given with a unit specified explicitly, for example `'120 ms'` for a time value, and they will be converted to whatever the parameter's actual unit is. Note that the value must be written as a string (with quotes) to use this feature. The unit name is case-sensitive, and there can be white-space between the numeric value and the unit.
 - Valid memory units are `B` (bytes), `kB` (kilobytes), `MB` (megabytes), `GB` (gigabytes), and `TB` (terabytes). The multiplier for memory units is 1024, not 1000.
 - Valid time units are `us` (microseconds), `ms` (milliseconds), `s` (seconds), `min` (minutes), `h` (hours), and `d` (days).

If a fractional value is specified with a unit, it will be rounded to a multiple of the next smaller unit if there is one. For example, `30.1 GB` will be converted to `30822 MB` not `32319628902 B`. If the parameter is of integer type, a final rounding to integer occurs after any unit conversion.

- *Enumerated*: Enumerated-type parameters are written in the same way as string parameters, but are restricted to have one of a limited set of values. The values allowable for such a parameter can be found from `pg_settings.enumvals`. Enum parameter values are case-insensitive.

19.1.2. Parameter Interaction via the Configuration File

The most fundamental way to set these parameters is to edit the file `postgresql.conf`, which is normally kept in the data directory, and `sdmspec.json` for a distributed system's configuration. A default copy is installed when the database cluster directory is initialized. An example of what this file might look like is:

```
# This is a comment
log_connections = yes
log_destination = 'syslog'
search_path = '$user', public
shared_buffers = 128MB
```

One parameter is specified per line. The equal sign between name and value is optional. Whitespace is insignificant (except within a quoted parameter value) and blank lines are ignored. Hash marks (#) designate the remainder of the line as a comment. Parameter values that are not simple identifiers or numbers must be single-quoted. To embed a single quote in a parameter value, write either two quotes (preferred) or backslash-quote. If the file contains multiple entries for the same parameter, all but the last one are ignored.

Parameters set in this way provide default values for the cluster. The settings seen by active sessions will be these values unless they are overridden. The following sections describe ways in which the administrator or user can override these defaults.

The configuration file is reread whenever the main server process receives a SIGHUP signal; this signal is most easily sent by running `pg_ctl reload` from the command line or by calling the SQL function `pg_reload_conf()`. The main server process also propagates this signal to all currently running server processes, so that existing sessions also adopt the new values (this will happen after they complete any currently-executing client command). Alternatively, you can send the signal to a single server process directly. Some parameters can only be set at server start; any changes to their entries in the configuration file will be ignored until the server is restarted. Invalid parameter settings in the configuration file are likewise ignored (but logged) during SIGHUP processing.

In addition to `postgresql.conf`, a Postgres Pro data directory contains a file `postgresql.auto.conf`, which has the same format as `postgresql.conf` but is intended to be edited automatically, not manually. Settings in `postgresql.auto.conf` override those in `postgresql.conf`.

External tools may also modify `postgresql.auto.conf`. It is not recommended to do this while the server is running unless `allow_alter_system` is set to `off`, since a concurrent `ALTER SYSTEM` command could overwrite such changes. Such tools might simply append new settings to the end, or they might choose to remove duplicate settings and/or comments (as `ALTER SYSTEM` will).

The system view `pg_file_settings` can be helpful for pre-testing changes to the configuration files, or for diagnosing problems if a SIGHUP signal did not have the desired effects.

19.1.3. Parameter Interaction via SQL

Postgres Pro provides three SQL commands to establish configuration defaults. The already-mentioned `ALTER SYSTEM` command provides an SQL-accessible means of changing global defaults; it is functionally equivalent to editing `postgresql.conf`. In addition, there are two commands that allow setting of defaults on a per-database or per-role basis:

- The `ALTER DATABASE` command allows global settings to be overridden on a per-database basis.
- The `ALTER ROLE` command allows both global and per-database settings to be overridden with user-specific values.

Values set with `ALTER DATABASE` and `ALTER ROLE` are applied only when starting a fresh database session. They override values obtained from the configuration files or server command line, and constitute defaults for the rest of the session. Note that some settings cannot be changed after server start, and so cannot be set with these commands (or the ones listed below).

Once a client is connected to the database, Postgres Pro provides two additional SQL commands (and equivalent functions) to interact with session-local configuration settings:

- The `SHOW` command allows inspection of the current value of any parameter. The corresponding SQL function is `current_setting(setting_name text)` (see [Section 9.28.1](#)).
- The `SET` command allows modification of the current value of those parameters that can be set locally to a session; it has no effect on other sessions. Many parameters can be set this way by any user, but some can only be set by superusers and users who have been granted `SET` privilege on that parameter. The corresponding SQL function is `set_config(setting_name, new_value, is_local)` (see [Section 9.28.1](#)).

In addition, the system view `pg_settings` can be used to view and change session-local values:

- Querying this view is similar to using `SHOW ALL` but provides more detail. It is also more flexible, since it's possible to specify filter conditions or join against other relations.
- Using `UPDATE` on this view, specifically updating the `setting` column, is the equivalent of issuing `SET` commands. For example, the equivalent of

```
SET configuration_parameter TO DEFAULT;
```

is:

```
UPDATE pg_settings SET setting = reset_val WHERE name = 'configuration_parameter';
```

19.1.4. Parameter Interaction via the Shell

In addition to setting global defaults or attaching overrides at the database or role level, you can pass settings to Postgres Pro via shell facilities. Both the server and libpq client library accept parameter values via the shell.

- During server startup, parameter settings can be passed to the `postgres` command via the `-c name=value` command-line parameter, or its equivalent `--name=value` variation. For example,

```
postgres -c log_connections=yes --log-destination='syslog'
```

Settings provided in this way override those set via `postgresql.conf` or `ALTER SYSTEM`, so they cannot be changed globally without restarting the server.

- When starting a client session via libpq, parameter settings can be specified using the `PGOPTIONS` environment variable. Settings established in this way constitute defaults for the life of the session, but do not affect other sessions. For historical reasons, the format of `PGOPTIONS` is similar to that used when launching the `postgres` command; specifically, the `-c`, or prepended `--`, before the name must be specified. For example,

```
env PGOPTIONS="-c geqo=off --statement-timeout=5min" psql
```

Other clients and libraries might provide their own mechanisms, via the shell or otherwise, that allow the user to alter session settings without direct use of SQL commands.

19.1.5. Managing `postgresql.conf` Configuration File Contents

Postgres Pro provides several features for breaking down complex `postgresql.conf` files into sub-files. These features are especially useful when managing multiple servers with related, but not identical, configurations. Note that you cannot edit it manually.

In addition to individual parameter settings, the `postgresql.conf` file can contain *include directives*, which specify another file to read and process as if it were inserted into the configuration file at this point. This feature allows a configuration file to be divided into physically separate parts. Include directives simply look like:

```
include 'filename'
```

If the file name is not an absolute path, it is taken as relative to the directory containing the referencing configuration file. Inclusions can be nested.

There is also an `include_if_exists` directive, which acts the same as the `include` directive, except when the referenced file does not exist or cannot be read. A regular `include` will consider this an error condition, but `include_if_exists` merely logs a message and continues processing the referencing configuration file.

The `postgresql.conf` file can also contain `include_dir` directives, which specify an entire directory of configuration files to include. These look like

```
include_dir 'directory'
```

Non-absolute directory names are taken as relative to the directory containing the referencing configuration file. Within the specified directory, only non-directory files whose names end with the suffix `.conf` will be included. File names that start with the `.` character are also ignored, to prevent mistakes

since such files are hidden on some platforms. Multiple files within an include directory are processed in file name order (according to C locale rules, i.e., numbers before letters, and uppercase letters before lowercase ones).

Include files or directories can be used to logically separate portions of the database configuration, rather than having a single large `postgresql.conf` file. Consider a company that has two database servers, each with a different amount of memory. There are likely elements of the configuration both will share, for things such as logging. But memory-related parameters on the server will vary between the two. And there might be server specific customizations, too. One way to manage this situation is to break the custom configuration changes for your site into three files. You could add this to the end of your `postgresql.conf` file to include them:

```
include 'shared.conf'
include 'memory.conf'
include 'server.conf'
```

All systems would have the same `shared.conf`. Each server with a particular amount of memory could share the same `memory.conf`; you might have one for all servers with 8GB of RAM, another for those having 16GB. And finally `server.conf` could have truly server-specific configuration information in it.

Another possibility is to create a configuration file directory and put this information into files there. For example, a `conf.d` directory could be referenced at the end of `postgresql.conf`:

```
include_dir 'conf.d'
```

Then you could name the files in the `conf.d` directory like this:

```
00shared.conf
01memory.conf
02server.conf
```

This naming convention establishes a clear order in which these files will be loaded. This is important because only the last setting encountered for a particular parameter while the server is reading configuration files will be used. In this example, something set in `conf.d/02server.conf` would override a value set in `conf.d/01memory.conf`.

You might instead use this approach to naming the files descriptively:

```
00shared.conf
01memory-8GB.conf
02server-foo.conf
```

This sort of arrangement gives a unique name for each configuration file variation. This can help eliminate ambiguity when several servers have their configurations all stored in one place, such as in a version control repository. (Storing database configuration files under version control is another good practice to consider.)

19.2. File Locations

In addition to the `postgresql.conf` file already mentioned, Postgres Pro uses two other manually-edited configuration files, which control client authentication (their use is discussed in [Chapter 20](#)). By default, all three configuration files are stored in the database cluster's data directory. The parameters described in this section allow the configuration files to be placed elsewhere. (Doing so can ease administration. In particular it is often easier to ensure that the configuration files are properly backed-up when they are kept separate.)

`data_directory (string)`

Specifies the directory to use for data storage. This parameter can only be set at server start.

`config_file (string)`

Specifies the main server configuration file (customarily called `postgresql.conf`). This parameter can only be set on the `postgres` command line.

`hba_file (string)`

Specifies the configuration file for host-based authentication (customarily called `pg_hba.conf`). This parameter can only be set at server start.

`ident_file (string)`

Specifies the configuration file for user name mapping (customarily called `pg_ident.conf`). This parameter can only be set at server start. See also [Section 20.2](#).

`external_pid_file (string)`

Specifies the name of an additional process-ID (PID) file that the server should create for use by server administration programs. This parameter can only be set at server start.

In a default installation, none of the above parameters are set explicitly. Instead, the data directory is specified by the `-D` command-line option or the `PGDATA` environment variable, and the configuration files are all found within the data directory.

If you wish to keep the configuration files elsewhere than the data directory, the `postgres -D` command-line option or `PGDATA` environment variable must point to the directory containing the configuration files, and the `data_directory` parameter must be set in `postgresql.conf` (or on the command line) to show where the data directory is actually located. Notice that `data_directory` overrides `-D` and `PGDATA` for the location of the data directory, but not for the location of the configuration files.

If you wish, you can specify the configuration file names and locations individually using the parameters `config_file`, `hba_file` and/or `ident_file`. `config_file` can only be specified on the `postgres` command line, but the others can be set within the main configuration file. If all three parameters plus `data_directory` are explicitly set, then it is not necessary to specify `-D` or `PGDATA`.

When setting any of these parameters, a relative path will be interpreted with respect to the directory in which `postgres` is started.

19.3. Connections and Authentication

19.3.1. Connection Settings

`listen_addresses (string)`

Specifies the TCP/IP address(es) on which the server is to listen for connections from client applications. The value takes the form of a comma-separated list of host names and/or numeric IP addresses. The special entry `*` corresponds to all available IP interfaces. The entry `0.0.0.0` allows listening for all IPv4 addresses and `::` allows listening for all IPv6 addresses. If the list is empty, the server does not listen on any IP interface at all, in which case only Unix-domain sockets can be used to connect to it. If the list is not empty, the server will start if it can listen on at least one TCP/IP address. A warning will be emitted for any TCP/IP address which cannot be opened. The default value is `localhost`, which allows only local TCP/IP “loopback” connections to be made.

While client authentication ([Chapter 20](#)) allows fine-grained control over who can access the server, `listen_addresses` controls which interfaces accept connection attempts, which can help prevent repeated malicious connection requests on insecure network interfaces. This parameter can only be set at server start.

`max_connections (integer)`

Determines the maximum number of concurrent connections to the database server. The default is typically 100 connections, but might be less if your kernel settings will not support it (as determined during `initdb`). This parameter can only be set at server start.

Postgres Pro sizes certain resources based directly on the value of `max_connections`. Increasing its value leads to higher allocation of those resources, including shared memory.

When running a standby server, you must set this parameter to the same or higher value than on the primary server. Otherwise, queries will not be allowed in the standby server.

`reserved_connections` (integer)

Determines the number of connection “slots” that are reserved for connections by roles with privileges of the `pg_use_reserved_connections` role. Whenever the number of free connection slots is greater than `superuser_reserved_connections` but less than or equal to the sum of `superuser_reserved_connections` and `reserved_connections`, new connections will be accepted only for superusers and roles with privileges of `pg_use_reserved_connections`. If `superuser_reserved_connections` or fewer connection slots are available, new connections will be accepted only for superusers.

The default value is zero connections. The value must be less than `max_connections` minus `superuser_reserved_connections`. This parameter can only be set at server start.

`superuser_reserved_connections` (integer)

Determines the number of connection “slots” that are reserved for connections by Postgres Pro superusers. At most `max_connections` connections can ever be active simultaneously. Whenever the number of active concurrent connections is at least `max_connections` minus `superuser_reserved_connections`, new connections will be accepted only for superusers. The connection slots reserved by this parameter are intended as final reserve for emergency use after the slots reserved by `reserved_connections` have been exhausted.

The default value is three connections. The value must be less than `max_connections` minus `reserved_connections`. This parameter can only be set at server start.

`unix_socket_directories` (string)

Specifies the directory of the Unix-domain socket(s) on which the server is to listen for connections from client applications. Multiple sockets can be created by listing multiple directories separated by commas. Whitespace between entries is ignored; surround a directory name with double quotes if you need to include whitespace or commas in the name. An empty value specifies not listening on any Unix-domain sockets, in which case only TCP/IP sockets can be used to connect to the server.

A value that starts with `@` specifies that a Unix-domain socket in the abstract namespace should be created (currently supported on Linux only). In that case, this value does not specify a “directory” but a prefix from which the actual socket name is computed in the same manner as for the file-system namespace. While the abstract socket name prefix can be chosen freely, since it is not a file-system location, the convention is to nonetheless use file-system-like values such as `@/tmp`.

The default value is normally `/tmp`, but that can be changed at build time. On Windows, the default is empty, which means no Unix-domain socket is created by default. This parameter can only be set at server start.

In addition to the socket file itself, which is named `.s.PGSQL.nnnn` where `nnnn` is the server's port number, an ordinary file named `.s.PGSQL.nnnn.lock` will be created in each of the `unix_socket_directories` directories. Neither file should ever be removed manually. For sockets in the abstract namespace, no lock file is created.

`unix_socket_group` (string)

Sets the owning group of the Unix-domain socket(s). (The owning user of the sockets is always the user that starts the server.) In combination with the parameter `unix_socket_permissions` this can be used as an additional access control mechanism for Unix-domain connections. By default this is the empty string, which uses the default group of the server user. This parameter can only be set at server start.

This parameter is not supported on Windows. Any setting will be ignored. Also, sockets in the abstract namespace have no file owner, so this setting is also ignored in that case.

`unix_socket_permissions` (integer)

Sets the access permissions of the Unix-domain socket(s). Unix-domain sockets use the usual Unix file system permission set. The parameter value is expected to be a numeric mode specified in the format accepted by the `chmod` and `umask` system calls. (To use the customary octal format the number must start with a 0 (zero).)

The default permissions are `0777`, meaning anyone can connect. Reasonable alternatives are `0770` (only user and group, see also `unix_socket_group`) and `0700` (only user). (Note that for a Unix-domain socket, only write permission matters, so there is no point in setting or revoking read or execute permissions.)

This access control mechanism is independent of the one described in [Chapter 20](#).

This parameter can only be set at server start.

This parameter is irrelevant on systems, notably Solaris as of Solaris 10, that ignore socket permissions entirely. There, one can achieve a similar effect by pointing `unix_socket_directories` to a directory having search permission limited to the desired audience.

Sockets in the abstract namespace have no file permissions, so this setting is also ignored in that case.

`bonjour` (boolean)

Enables advertising the server's existence via Bonjour. The default is off. This parameter can only be set at server start.

`bonjour_name` (string)

Specifies the Bonjour service name. The computer name is used if this parameter is set to the empty string `' '` (which is the default). This parameter is ignored if the server was not compiled with Bonjour support. This parameter can only be set at server start.

19.3.2. TCP Settings

`tcp_keepalives_idle` (integer)

Specifies the amount of time with no network activity after which the operating system should send a TCP keepalive message to the client. If this value is specified without units, it is taken as seconds. A value of 0 (the default) selects the operating system's default. On Windows, setting a value of 0 will set this parameter to 2 hours, since Windows does not provide a way to read the system default value. This parameter is supported only on systems that support `TCP_KEEPIDL` or an equivalent socket option, and on Windows; on other systems, it must be zero. In sessions connected via a Unix-domain socket, this parameter is ignored and always reads as zero.

`tcp_keepalives_interval` (integer)

Specifies the amount of time after which a TCP keepalive message that has not been acknowledged by the client should be retransmitted. If this value is specified without units, it is taken as seconds. A value of 0 (the default) selects the operating system's default. On Windows, setting a value of 0 will set this parameter to 1 second, since Windows does not provide a way to read the system default value. This parameter is supported only on systems that support `TCP_KEEPINTVL` or an equivalent socket option, and on Windows; on other systems, it must be zero. In sessions connected via a Unix-domain socket, this parameter is ignored and always reads as zero.

`tcp_keepalives_count` (integer)

Specifies the number of TCP keepalive messages that can be lost before the server's connection to the client is considered dead. A value of 0 (the default) selects the operating system's default. This parameter is supported only on systems that support `TCP_KEEPCNT` or an equivalent socket option (which does not include Windows); on other systems, it must be zero. In sessions connected via a Unix-domain socket, this parameter is ignored and always reads as zero.

`tcp_user_timeout (integer)`

Specifies the amount of time that transmitted data may remain unacknowledged before the TCP connection is forcibly closed. If this value is specified without units, it is taken as milliseconds. A value of 0 (the default) selects the operating system's default. This parameter is supported only on systems that support `TCP_USER_TIMEOUT` (which does not include Windows); on other systems, it must be zero. In sessions connected via a Unix-domain socket, this parameter is ignored and always reads as zero.

`client_connection_check_interval (integer)`

Sets the time interval between optional checks that the client is still connected, while running queries. The check is performed by polling the socket, and allows long running queries to be aborted sooner if the kernel reports that the connection is closed.

This option relies on kernel events exposed by Linux, macOS, illumos and the BSD family of operating systems, and is not currently available on other systems.

If the value is specified without units, it is taken as milliseconds. The default value is 0, which disables connection checks. Without connection checks, the server will detect the loss of the connection only at the next interaction with the socket, when it waits for, receives or sends data.

For the kernel itself to detect lost TCP connections reliably and within a known timeframe in all scenarios including network failure, it may also be necessary to adjust the TCP keepalive settings of the operating system, or the [tcp_keepalives_idle](#), [tcp_keepalives_interval](#) and [tcp_keepalives_count](#) settings of Postgres Pro.

19.3.3. Authentication

`authentication_timeout (integer)`

Maximum amount of time allowed to complete client authentication. If a would-be client has not completed the authentication protocol in this much time, the server closes the connection. This prevents hung clients from occupying a connection indefinitely. If this value is specified without units, it is taken as seconds. The default is one minute (1m). This parameter can only be set in the `postgresql.conf` file or on the server command line.

`password_encryption (enum)`

When a password is specified in [CREATE ROLE](#) or [ALTER ROLE](#), this parameter determines the algorithm to use to encrypt the password. Possible values are `scram-sha-256`, which will encrypt the password with SCRAM-SHA-256, and `md5`, which stores the password as an MD5 hash. The default is `scram-sha-256`. Nevertheless, by default [pg-setup](#) sets `--auth-host` to `md5` but you can still pass another value in its `initdb` option, if necessary.

Note that older clients might lack support for the SCRAM authentication mechanism, and hence not work with passwords encrypted with SCRAM-SHA-256. See [Section 20.5](#) for more details.

`scram_iterations (integer)`

The number of computational iterations to be performed when encrypting a password using SCRAM-SHA-256. The default is 4096. A higher number of iterations provides additional protection against brute-force attacks on stored passwords, but makes authentication slower. Changing the value has no effect on existing passwords encrypted with SCRAM-SHA-256 as the iteration count is fixed at the time of encryption. In order to make use of a changed value, a new password must be set.

`krb_server_keyfile (string)`

Sets the location of the server's Kerberos key file. The default is `FILE:/usr/local/pgsql/etc/krb5.keytab` (where the directory part is whatever was specified as `sysconfdir` at build time; use `pg_config --sysconfdir` to determine that). If this parameter is set to an empty string, it is ignored and a system-dependent default is used. This parameter can only be set in the `postgresql.conf` file or on the server command line. See [Section 20.6](#) for more information.

`krb_caseins_users` (boolean)

Sets whether GSSAPI user names should be treated case-insensitively. The default is `off` (case sensitive). This parameter can only be set in the `postgresql.conf` file or on the server command line.

`gss_accept_delegation` (boolean)

Sets whether GSSAPI delegation should be accepted from the client. The default is `off` meaning credentials from the client will *not* be accepted. Changing this to `on` will make the server accept credentials delegated to it from the client. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.3.4. SSL

See [Section 18.9](#) for more information about setting up SSL. The configuration parameters for controlling transfer encryption using TLS protocols are named `ssl` for historic reasons, even though support for the SSL protocol has been deprecated. SSL is in this context used interchangeably with TLS.

`ssl` (boolean)

Enables SSL connections. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `off`.

`ssl_ca_file` (string)

Specifies the name of the file containing the SSL server certificate authority (CA). Relative paths are relative to the data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is empty, meaning no CA file is loaded, and client certificate verification is not performed.

`ssl_cert_file` (string)

Specifies the name of the file containing the SSL server certificate. Relative paths are relative to the data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `server.crt`.

`ssl_crl_file` (string)

Specifies the name of the file containing the SSL client certificate revocation list (CRL). Relative paths are relative to the data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is empty, meaning no CRL file is loaded (unless [ssl_crl_dir](#) is set).

`ssl_crl_dir` (string)

Specifies the name of the directory containing the SSL client certificate revocation list (CRL). Relative paths are relative to the data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is empty, meaning no CRLs are used (unless [ssl_crl_file](#) is set).

The directory needs to be prepared with the OpenSSL command `openssl rehash` or `c_rehash`. See its documentation for details.

When using this setting, CRLs in the specified directory are loaded on-demand at connection time. New CRLs can be added to the directory and will be used immediately. This is unlike [ssl_crl_file](#), which causes the CRL in the file to be loaded at server start time or when the configuration is reloaded. Both settings can be used together.

`ssl_key_file` (string)

Specifies the name of the file containing the SSL server private key. Relative paths are relative to the data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `server.key`.

ssl_ciphers (string)

Specifies a list of SSL cipher suites that are allowed to be used by SSL connections. See the ciphers manual page in the OpenSSL package for the syntax of this setting and a list of supported values. Only connections using TLS version 1.2 and lower are affected. There is currently no setting that controls the cipher choices used by TLS version 1.3 connections. The default value is `HIGH:MEDIUM:+3DES:!aNULL`. The default is usually a reasonable choice unless you have specific security requirements.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

Explanation of the default value:

HIGH

Cipher suites that use ciphers from HIGH group (e.g., AES, Camellia, 3DES)

MEDIUM

Cipher suites that use ciphers from MEDIUM group (e.g., RC4, SEED)

+3DES

The OpenSSL default order for HIGH is problematic because it orders 3DES higher than AES128. This is wrong because 3DES offers less security than AES128, and it is also much slower. +3DES reorders it after all other HIGH and MEDIUM ciphers.

!aNULL

Disables anonymous cipher suites that do no authentication. Such cipher suites are vulnerable to MITM attacks and therefore should not be used.

Available cipher suite details will vary across OpenSSL versions. Use the command `openssl ciphers -v 'HIGH:MEDIUM:+3DES:!aNULL'` to see actual details for the currently installed OpenSSL version. Note that this list is filtered at run time based on the server key type.

ssl_prefer_server_ciphers (boolean)

Specifies whether to use the server's SSL cipher preferences, rather than the client's. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `on`.

Postgres Pro versions before 9.4 do not have this setting and always use the client's preferences. This setting is mainly for backward compatibility with those versions. Using the server's preferences is usually better because it is more likely that the server is appropriately configured.

ssl_ecdh_curve (string)

Specifies the name of the curve to use in ECDH key exchange. It needs to be supported by all clients that connect. It does not need to be the same curve used by the server's Elliptic Curve key. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `prime256v1`.

OpenSSL names for the most common curves are: `prime256v1` (NIST P-256), `secp384r1` (NIST P-384), `secp521r1` (NIST P-521). The full list of available curves can be shown with the command `openssl ecparam -list_curves`. Not all of them are usable in TLS though.

ssl_min_protocol_version (enum)

Sets the minimum SSL/TLS protocol version to use. Valid values are currently: `TLSv1`, `TLSv1.1`, `TLSv1.2`, `TLSv1.3`. Older versions of the OpenSSL library do not support all values; an error will be raised if an unsupported setting is chosen. Protocol versions before TLS 1.0, namely SSL version 2 and 3, are always disabled.

The default is `TLSv1.2`, which satisfies industry best practices as of this writing.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`ssl_max_protocol_version` (enum)

Sets the maximum SSL/TLS protocol version to use. Valid values are as for [ssl_min_protocol_version](#), with addition of an empty string, which allows any protocol version. The default is to allow any version. Setting the maximum protocol version is mainly useful for testing or if some component has issues working with a newer protocol.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`ssl_dh_params_file` (string)

Specifies the name of the file containing Diffie-Hellman parameters used for so-called ephemeral DH family of SSL ciphers. The default is empty, in which case compiled-in default DH parameters used. Using custom DH parameters reduces the exposure if an attacker manages to crack the well-known compiled-in DH parameters. You can create your own DH parameters file with the command `openssl dhparam -out dhparams.pem 2048`.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`ssl_passphrase_command` (string)

Sets an external command to be invoked when a passphrase for decrypting an SSL file such as a private key needs to be obtained. By default, this parameter is empty, which means the built-in prompting mechanism is used.

The command must print the passphrase to the standard output and exit with code 0. In the parameter value, `%p` is replaced by a prompt string. (Write `%%` for a literal `%`.) Note that the prompt string will probably contain whitespace, so be sure to quote adequately. A single newline is stripped from the end of the output if present.

The command does not actually have to prompt the user for a passphrase. It can read it from a file, obtain it from a keychain facility, or similar. It is up to the user to make sure the chosen mechanism is adequately secure.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`ssl_passphrase_command_supports_reload` (boolean)

This parameter determines whether the passphrase command set by `ssl_passphrase_command` will also be called during a configuration reload if a key file needs a passphrase. If this parameter is off (the default), then `ssl_passphrase_command` will be ignored during a reload and the SSL configuration will not be reloaded if a passphrase is needed. That setting is appropriate for a command that requires a TTY for prompting, which might not be available when the server is running. Setting this parameter to on might be appropriate if the passphrase is obtained from a file, for example.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.4. Resource Consumption

19.4.1. Memory

`shared_buffers` (integer)

Sets the amount of memory the database server uses for shared memory buffers. The default is typically 128 megabytes (128MB), but might be less if your kernel settings will not support it (as determined during `initdb`). This setting must be at least 128 kilobytes. However, settings significantly higher than the minimum are usually needed for good performance. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. (Non-default values of `BLCKSZ` change the minimum value.) This parameter can only be set at server start.

If you have a dedicated database server with 1GB or more of RAM, a reasonable starting value for `shared_buffers` is 25% of the memory in your system. There are some workloads where even larger

settings for `shared_buffers` are effective, but because Postgres Pro also relies on the operating system cache, it is unlikely that an allocation of more than 40% of RAM to `shared_buffers` will work better than a smaller amount. Larger settings for `shared_buffers` usually require a corresponding increase in `max_wal_size`, in order to spread out the process of writing large quantities of new or changed data over a longer period of time.

On systems with less than 1GB of RAM, a smaller percentage of RAM is appropriate, so as to leave adequate space for the operating system.

`huge_pages` (enum)

Controls whether huge pages are requested for the main shared memory area. Valid values are `try` (the default), `on`, and `off`. With `huge_pages` set to `try`, the server will try to request huge pages, but fall back to the default if that fails. With `on`, failure to request huge pages will prevent the server from starting up. With `off`, huge pages will not be requested. The actual state of huge pages is indicated by the server variable [huge_pages_status](#).

At present, this setting is supported only on Linux and Windows. The setting is ignored on other systems when set to `try`. On Linux, it is only supported when `shared_memory_type` is set to `mmap` (the default).

The use of huge pages results in smaller page tables and less CPU time spent on memory management, increasing performance. For more details about using huge pages on Linux, see [Section 18.4.5](#).

Huge pages are known as large pages on Windows. To use them, you need to assign the user right “Lock pages in memory” to the Windows user account that runs Postgres Pro. You can use Windows Group Policy tool (`gpedit.msc`) to assign the user right “Lock pages in memory”. To start the database server on the command prompt as a standalone process, not as a Windows service, the command prompt must be run as an administrator or User Access Control (UAC) must be disabled. When the UAC is enabled, the normal command prompt revokes the user right “Lock pages in memory” when started.

Note that this setting only affects the main shared memory area. Operating systems such as Linux, FreeBSD, and Illumos can also use huge pages (also known as “super” pages or “large” pages) automatically for normal memory allocation, without an explicit request from Postgres Pro. On Linux, this is called “transparent huge pages”(THP). That feature has been known to cause performance degradation with Postgres Pro for some users on some Linux versions, so its use is currently discouraged (unlike explicit use of `huge_pages`).

`xml_parse_huge` (boolean)

Allows allocating more memory for processing XML data. By default, the maximum amount of memory allocated per node is 10 MB. With `xml_parse_huge` enabled, this limit is increased to 1 GB.

Default: `off`

`huge_page_size` (integer)

Controls the size of huge pages, when they are enabled with [huge_pages](#). The default is zero (0). When set to 0, the default huge page size on the system will be used. This parameter can only be set at server start.

Some commonly available page sizes on modern 64 bit server architectures include: 2MB and 1GB (Intel and AMD), 16MB and 16GB (IBM POWER), and 64kB, 2MB, 32MB and 1GB (ARM). For more information about usage and support, see [Section 18.4.5](#).

Non-default settings are currently supported only on Linux.

`temp_buffers` (integer)

Sets the maximum amount of memory used for temporary buffers within each database session. These are session-local buffers used only for access to temporary tables. If this value is specified without

units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default is eight megabytes (8MB). (If `BLCKSZ` is not 8kB, the default value scales proportionally to it.) This setting can be changed within individual sessions, but only before the first use of temporary tables within the session; subsequent attempts to change the value will have no effect on that session.

A session will allocate temporary buffers as needed up to the limit given by `temp_buffers`. The cost of setting a large value in sessions that do not actually need many temporary buffers is only a buffer descriptor, or about 64 bytes, per increment in `temp_buffers`. However if a buffer is actually used an additional 8192 bytes will be consumed for it (or in general, `BLCKSZ` bytes).

`max_prepared_transactions` (integer)

Sets the maximum number of transactions that can be in the “prepared” state simultaneously (see [PREPARE TRANSACTION](#)). Setting this parameter to zero (which is the default) disables the prepared-transaction feature. This parameter can only be set at server start.

For distributed system specifics, see [max_prepared_transactions](#).

If you are not planning to use prepared transactions, this parameter should be set to zero to prevent accidental creation of prepared transactions. If you are using prepared transactions, you will probably want `max_prepared_transactions` to be at least as large as [max_connections](#), so that every session can have a prepared transaction pending.

When running a standby server, you must set this parameter to the same or higher value than on the primary server. Otherwise, queries will not be allowed in the standby server.

`max_autonomous_transactions` (integer)

Sets the maximum number of autonomous transactions that can be used simultaneously in all sessions of a Postgres Pro instance. If this value is exceeded, the error is raised that an autonomous transaction cannot be started.

Default: 100

`work_mem` (integer)

Sets the base maximum amount of memory to be used by a query operation (such as a sort or hash table) before writing to temporary disk files. If this value is specified without units, it is taken as kilobytes. The default value is four megabytes (4MB). Note that a complex query might perform several sort and hash operations at the same time, with each operation generally being allowed to use as much memory as this value specifies before it starts to write data into temporary files. Also, several running sessions could be doing such operations concurrently. Therefore, the total memory used could be many times the value of `work_mem`; it is necessary to keep this fact in mind when choosing the value. Sort operations are used for `ORDER BY`, `DISTINCT`, and merge joins. Hash tables are used in hash joins, hash-based aggregation, memoize nodes and hash-based processing of `IN` subqueries.

Hash-based operations are generally more sensitive to memory availability than equivalent sort-based operations. The memory limit for a hash table is computed by multiplying `work_mem` by `hash_mem_multiplier`. This makes it possible for hash-based operations to use an amount of memory that exceeds the usual `work_mem` base amount.

Note

The Postgres Pro planner may incorrectly estimate the hash table size when doing hash-based operations. If the hash table is organized in a such manner that one or several batches are so big that they do not fit into the allocated memory, exceeding the memory limit, and the hash table can not be splitted into more batches, then the planner will try to place each batch into the memory in order to complete the query even if a single batch exceeds the aforementioned limit.

Therefore, the memory consumption could dramatically increase when doing hash-based operations on large datasets that produce hash tables, which can not be splitted into a large number of batches.

`hash_mem_multiplier` (floating point)

Used to compute the maximum amount of memory that hash-based operations can use. The final limit is determined by multiplying `work_mem` by `hash_mem_multiplier`. The default value is 2.0, which makes hash-based operations use twice the usual `work_mem` base amount.

Consider increasing `hash_mem_multiplier` in environments where spilling by query operations is a regular occurrence, especially when simply increasing `work_mem` results in memory pressure (memory pressure typically takes the form of intermittent out of memory errors). The default setting of 2.0 is often effective with mixed workloads. Higher settings in the range of 2.0 - 8.0 or more may be effective in environments where `work_mem` has already been increased to 40MB or more.

`maintenance_work_mem` (integer)

Specifies the maximum amount of memory to be used by maintenance operations, such as `VACUUM`, `CREATE INDEX`, and `ALTER TABLE ADD FOREIGN KEY`. If this value is specified without units, it is taken as kilobytes. It defaults to 64 megabytes (64MB). Since only one of these operations can be executed at a time by a database session, and an installation normally doesn't have many of them running concurrently, it's safe to set this value significantly larger than `work_mem`. Larger settings might improve performance for vacuuming and for restoring database dumps.

Note that when autovacuum runs, up to `autovacuum_max_workers` times this memory may be allocated, so be careful not to set the default value too high. It may be useful to control for this by separately setting `autovacuum_work_mem`.

`autovacuum_work_mem` (integer)

Specifies the maximum amount of memory to be used by each autovacuum worker process. If this value is specified without units, it is taken as kilobytes. It defaults to -1, indicating that the value of `maintenance_work_mem` should be used instead. The setting has no effect on the behavior of `VACUUM` when run in other contexts. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`vacuum_buffer_usage_limit` (integer)

Specifies the size of the *Buffer Access Strategy* used by the `VACUUM` and `ANALYZE` commands. A setting of 0 will allow the operation to use any number of `shared_buffers`. Otherwise valid sizes range from 128 kB to 16 GB. If the specified size would exceed 1/8 the size of `shared_buffers`, the size is silently capped to that value. The default value is 2MB. If this value is specified without units, it is taken as kilobytes. This parameter can be set at any time. It can be overridden for `VACUUM` and `ANALYZE` when passing the `BUFFER_USAGE_LIMIT` option. Higher settings can allow `VACUUM` and `ANALYZE` to run more quickly, but having too large a setting may cause too many other useful pages to be evicted from shared buffers.

`logical_decoding_work_mem` (integer)

Specifies the maximum amount of memory to be used by logical decoding, before some of the decoded changes are written to local disk. This limits the amount of memory used by logical streaming replication connections. It defaults to 64 megabytes (64MB). Since each replication connection only uses a single buffer of this size, and an installation normally doesn't have many such connections concurrently (as limited by `max_wal_senders`), it's safe to set this value significantly higher than `work_mem`, reducing the amount of decoded changes written to disk.

`commit_timestamp_buffers` (integer)

Specifies the amount of memory to use to cache the contents of `pg_commit_ts` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The

default value is 0, which requests `shared_buffers/512` up to 1024 blocks, but not fewer than 16 blocks. This parameter can only be set at server start.

`multixact_member_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_multixact/members` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 32. This parameter can only be set at server start.

`multixact_offset_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_multixact/offsets` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 16. This parameter can only be set at server start.

`notify_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_notify` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 16. This parameter can only be set at server start.

`serializable_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_serial` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 32. This parameter can only be set at server start.

`subtransaction_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_subtrans` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 0, which requests `shared_buffers/512` up to 1024 blocks, but not fewer than 16 blocks. This parameter can only be set at server start.

`transaction_buffers` (integer)

Specifies the amount of shared memory to use to cache the contents of `pg_xact` (see [Table 66.1](#)). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default value is 0, which requests `shared_buffers/512` up to 1024 blocks, but not fewer than 16 blocks. This parameter can only be set at server start.

`max_stack_depth` (integer)

Specifies the maximum safe depth of the server's execution stack. The ideal setting for this parameter is the actual stack size limit enforced by the kernel (as set by `ulimit -s` or local equivalent), less a safety margin of a megabyte or so. The safety margin is needed because the stack depth is not checked in every routine in the server, but only in key potentially-recursive routines. If this value is specified without units, it is taken as kilobytes. The default setting is two megabytes (2MB), which is conservatively small and unlikely to risk crashes. However, it might be too small to allow execution of complex functions. Only superusers and users with the appropriate `SET` privilege can change this setting.

Setting `max_stack_depth` higher than the actual kernel limit will mean that a runaway recursive function can crash an individual backend process. On platforms where Postgres Pro can determine the kernel limit, the server will not allow this variable to be set to an unsafe value. However, not all platforms provide the information, so caution is recommended in selecting a value.

`shared_memory_type` (enum)

Specifies the shared memory implementation that the server should use for the main shared memory region that holds Postgres Pro's shared buffers and other shared data. Possible values are `mmap` (for anonymous shared memory allocated using `mmap`), `sysv` (for System V shared memory allocated via `shmget`) and `windows` (for Windows shared memory). Not all values are supported on all platforms; the first supported option is the default for that platform. The use of the `sysv` option, which is not

the default on any platform, is generally discouraged because it typically requires non-default kernel settings to allow for large allocations (see [Section 18.4.1](#)).

`dynamic_shared_memory_type` (enum)

Specifies the dynamic shared memory implementation that the server should use. Possible values are `posix` (for POSIX shared memory allocated using `shm_open`), `sysv` (for System V shared memory allocated via `shmget`), `windows` (for Windows shared memory), and `mmap` (to simulate shared memory using memory-mapped files stored in the data directory). Not all values are supported on all platforms; the first supported option is usually the default for that platform. The use of the `mmap` option, which is not the default on any platform, is generally discouraged because the operating system may write modified pages back to disk repeatedly, increasing system I/O load; however, it may be useful for debugging, when the `pg_dynshmem` directory is stored on a RAM disk, or when other shared memory facilities are not available.

`min_dynamic_shared_memory` (integer)

Specifies the amount of memory that should be allocated at server startup for use by parallel queries. When this memory region is insufficient or exhausted by concurrent queries, new parallel queries try to allocate extra shared memory temporarily from the operating system using the method configured with `dynamic_shared_memory_type`, which may be slower due to memory management overheads. Memory that is allocated at startup with `min_dynamic_shared_memory` is affected by the `huge_pages` setting on operating systems where that is supported, and may be more likely to benefit from larger pages on operating systems where that is managed automatically. The default value is 0 (none). This parameter can only be set at server start.

`enable_temp_memory_catalog` (boolean)

Enables or disables the use of an in-memory system catalog for temporary tables. The default value is `off`. This feature is useful to reduce an overhead of extra records in the catalog related to temporary tables. Thus, it improves system performance in case of massive use of temporary tables. Note that this feature is currently experimental.

Set the `enable_temp_memory_catalog` parameter to `on` only when the [enable_parallel_tempables](#) parameter is disabled. The corresponding features are mutually exclusive.

19.4.2. Disk

`temp_file_limit` (integer)

Specifies the maximum amount of disk space that a process can use for temporary files, such as sort and hash temporary files, or the storage file for a held cursor. A transaction attempting to exceed this limit will be canceled. If this value is specified without units, it is taken as kilobytes. `-1` (the default) means no limit. Only superusers and users with the appropriate `SET` privilege can change this setting.

This setting constrains the total space used at any instant by all temporary files used by a given Postgres Pro process. It should be noted that disk space used for explicit temporary tables, as opposed to temporary files used behind-the-scenes in query execution, does *not* count against this limit.

`max_notify_queue_pages` (integer)

Specifies the maximum amount of allocated pages for [NOTIFY](#) / [LISTEN](#) queue. The default value is 1048576. For 8 KB pages it allows to consume up to 8 GB of disk space.

`plan_cache_lru_memsizes` (integer)

Specifies the maximum amount of memory that can be used by plans of prepared statements. Having such a limit is useful for sessions with multiple prepared statements that may otherwise use too much memory. If this limit is reached, Postgres Pro Shardman evicts the generic plans of the least recently used statements from memory, keeping only the corresponding parsed queries. If such a statement is called again later, it has to be analyzed and planned again. To keep plans in memory for all prepared statements as long as possible, set this parameter to 0. If the value for this parameter is specified without units, it is taken as kilobytes. It is worth noting that the total amount of memory used by

the saved queries may be larger than the value of this parameter because it takes into account only the generated plans.

If this parameter is set together with [plan_cache_lru_size](#), the first reached limit takes effect.

Default: 8MB

`plan_cache_lru_size` (integer)

Same as [plan_cache_lru_memsize](#), but limits the number of prepared statements instead of the memory size. Zero means no limit.

If this parameter is set together with [plan_cache_lru_memsize](#), the first reached limit takes effect.

Default: 0

`max_backend_memory` (integer)

Specifies the maximum amount of memory that can be allocated to a backend. Zero means no limit.

Default: 0

`temp_table_max_size` (integer)

Specifies the maximum amount of disk space that a single temporary table can occupy. If this value is specified without units, it is taken as megabytes. The value of zero means no limit. Only superusers and users with the appropriate `SET` privilege can change this setting.

19.4.3. Kernel Resource Usage

`max_files_per_process` (integer)

Sets the maximum number of simultaneously open files allowed to each server subprocess. The default is one thousand files. If the kernel is enforcing a safe per-process limit, you don't need to worry about this setting. But on some platforms (notably, most BSD systems), the kernel will allow individual processes to open many more files than the system can actually support if many processes all try to open that many files. If you find yourself seeing “Too many open files” failures, try reducing this setting. This parameter can only be set at server start.

19.4.4. Cost-based Vacuum Delay

During the execution of [VACUUM](#) and [ANALYZE](#) commands, the system maintains an internal counter that keeps track of the estimated cost of the various I/O operations that are performed. When the accumulated cost reaches a limit (specified by `vacuum_cost_limit`), the process performing the operation will sleep for a short period of time, as specified by `vacuum_cost_delay`. Then it will reset the counter and continue execution.

The intent of this feature is to allow administrators to reduce the I/O impact of these commands on concurrent database activity. There are many situations where it is not important that maintenance commands like `VACUUM` and `ANALYZE` finish quickly; however, it is usually very important that these commands do not significantly interfere with the ability of the system to perform other database operations. Cost-based vacuum delay provides a way for administrators to achieve this.

This feature is disabled by default for manually issued `VACUUM` commands. To enable it, set the `vacuum_cost_delay` variable to a nonzero value.

`vacuum_cost_delay` (floating point)

The amount of time that the process will sleep when the cost limit has been exceeded. If this value is specified without units, it is taken as milliseconds. The default value is zero, which disables the cost-based vacuum delay feature. Positive values enable cost-based vacuuming.

When using cost-based vacuuming, appropriate values for `vacuum_cost_delay` are usually quite small, perhaps less than 1 millisecond. While `vacuum_cost_delay` can be set to fractional-millisecond

values, such delays may not be measured accurately on older platforms. On such platforms, increasing `VACUUM`'s throttled resource consumption above what you get at 1ms will require changing the other vacuum cost parameters. You should, nonetheless, keep `vacuum_cost_delay` as small as your platform will consistently measure; large delays are not helpful.

`vacuum_cost_page_hit` (integer)

The estimated cost for vacuuming a buffer found in the shared buffer cache. It represents the cost to lock the buffer pool, lookup the shared hash table and scan the content of the page. The default value is one.

`vacuum_cost_page_miss` (integer)

The estimated cost for vacuuming a buffer that has to be read from disk. This represents the effort to lock the buffer pool, lookup the shared hash table, read the desired block in from the disk and scan its content. The default value is 2.

`vacuum_cost_page_dirty` (integer)

The estimated cost charged when vacuum modifies a block that was previously clean. It represents the extra I/O required to flush the dirty block out to disk again. The default value is 20.

`vacuum_cost_limit` (integer)

This is the accumulated cost that will cause the vacuuming process to sleep for `vacuum_cost_delay`. The default is 200.

Note

There are certain operations that hold critical locks and should therefore complete as quickly as possible. Cost-based vacuum delays do not occur during such operations. Therefore it is possible that the cost accumulates far higher than the specified limit. To avoid uselessly long delays in such cases, the actual delay is calculated as `vacuum_cost_delay * accumulated_balance / vacuum_cost_limit` with a maximum of `vacuum_cost_delay * 4`.

19.4.5. Background Writer

There is a separate server process called the *background writer*, whose function is to issue writes of “dirty” (new or modified) shared buffers. When the number of clean shared buffers appears to be insufficient, the background writer writes some dirty buffers to the file system and marks them as clean. This reduces the likelihood that server processes handling user queries will be unable to find clean buffers and have to write dirty buffers themselves. However, the background writer does cause a net overall increase in I/O load, because while a repeatedly-dirtied page might otherwise be written only once per checkpoint interval, the background writer might write it several times as it is dirtied in the same interval. The parameters discussed in this subsection can be used to tune the behavior for local needs.

`bgwriter_delay` (integer)

Specifies the delay between activity rounds for the background writer. In each round the writer issues writes for some number of dirty buffers (controllable by the following parameters). It then sleeps for the length of `bgwriter_delay`, and repeats. When there are no dirty buffers in the buffer pool, though, it goes into a longer sleep regardless of `bgwriter_delay`. If this value is specified without units, it is taken as milliseconds. The default value is 200 milliseconds (200ms). Note that on some systems, the effective resolution of sleep delays is 10 milliseconds; setting `bgwriter_delay` to a value that is not a multiple of 10 might have the same results as setting it to the next higher multiple of 10. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`bgwriter_lru_maxpages` (integer)

In each round, no more than this many buffers will be written by the background writer. Setting this to zero disables background writing. (Note that checkpoints, which are managed by a separate,

dedicated auxiliary process, are unaffected.) The default value is 100 buffers. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`bgwriter_lru_multiplier` (floating point)

The number of dirty buffers written in each round is based on the number of new buffers that have been needed by server processes during recent rounds. The average recent need is multiplied by `bgwriter_lru_multiplier` to arrive at an estimate of the number of buffers that will be needed during the next round. Dirty buffers are written until there are that many clean, reusable buffers available. (However, no more than `bgwriter_lru_maxpages` buffers will be written per round.) Thus, a setting of 1.0 represents a “just in time” policy of writing exactly the number of buffers predicted to be needed. Larger values provide some cushion against spikes in demand, while smaller values intentionally leave writes to be done by server processes. The default is 2.0. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`bgwriter_flush_after` (integer)

Whenever more than this amount of data has been written by the background writer, attempt to force the OS to issue these writes to the underlying storage. Doing so will limit the amount of dirty data in the kernel's page cache, reducing the likelihood of stalls when an `fsync` is issued at the end of a checkpoint, or when the OS writes data back in larger batches in the background. Often that will result in greatly reduced transaction latency, but there also are some cases, especially with workloads that are bigger than [shared_buffers](#), but smaller than the OS's page cache, where performance might degrade. This setting may have no effect on some platforms. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The valid range is between 0, which disables forced writeback, and 2MB. The default is 512kB on Linux, 0 elsewhere. (If `BLCKSZ` is not 8kB, the default and maximum values scale proportionally to it.) This parameter can only be set in the `postgresql.conf` file or on the server command line.

Smaller values of `bgwriter_lru_maxpages` and `bgwriter_lru_multiplier` reduce the extra I/O load caused by the background writer, but make it more likely that server processes will have to issue writes for themselves, delaying interactive queries.

19.4.6. Asynchronous Behavior

`backend_flush_after` (integer)

Whenever more than this amount of data has been written by a single backend, attempt to force the OS to issue these writes to the underlying storage. Doing so will limit the amount of dirty data in the kernel's page cache, reducing the likelihood of stalls when an `fsync` is issued at the end of a checkpoint, or when the OS writes data back in larger batches in the background. Often that will result in greatly reduced transaction latency, but there also are some cases, especially with workloads that are bigger than [shared_buffers](#), but smaller than the OS's page cache, where performance might degrade. This setting may have no effect on some platforms. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The valid range is between 0, which disables forced writeback, and 2MB. The default is 0, i.e., no forced writeback. (If `BLCKSZ` is not 8kB, the maximum value scales proportionally to it.)

`effective_io_concurrency` (integer)

Sets the number of concurrent disk I/O operations that Postgres Pro expects can be executed simultaneously. Raising this value will increase the number of I/O operations that any individual Postgres Pro session attempts to initiate in parallel. The allowed range is 1 to 1000, or zero to disable issuance of asynchronous I/O requests. Currently, this setting only affects bitmap heap scans.

For magnetic drives, a good starting point for this setting is the number of separate drives comprising a RAID 0 stripe or RAID 1 mirror being used for the database. (For RAID 5 the parity drive should not be counted.) However, if the database is often busy with multiple queries issued in concurrent sessions, lower values may be sufficient to keep the disk array busy. A value higher than needed to keep the disks busy will only result in extra CPU overhead. SSDs and other memory-based storage can often process many concurrent requests, so the best value might be in the hundreds.

Asynchronous I/O depends on an effective `posix_fadvise` function, which some operating systems lack. If the function is not present then setting this parameter to anything but zero will result in an error. On some operating systems (e.g., Solaris), the function is present but does not actually do anything.

The default is 1 on supported systems, otherwise 0. This value can be overridden for tables in a particular tablespace by setting the tablespace parameter of the same name (see [ALTER TABLESPACE](#)).

`maintenance_io_concurrency` (integer)

Similar to `effective_io_concurrency`, but used for maintenance work that is done on behalf of many client sessions.

The default is 10 on supported systems, otherwise 0. This value can be overridden for tables in a particular tablespace by setting the tablespace parameter of the same name (see [ALTER TABLESPACE](#)).

`io_combine_limit` (integer)

Controls the largest I/O size in operations that combine I/O. The default is 128kB.

`max_worker_processes` (integer)

Sets the maximum number of background processes that the cluster can support. This parameter can only be set at server start. The default is 8.

When running a standby server, you must set this parameter to the same or higher value than on the primary server. Otherwise, queries will not be allowed in the standby server.

When changing this value, consider also adjusting [max_parallel_workers](#), [max_parallel_maintenance_workers](#), and [max_parallel_workers_per_gather](#).

`max_parallel_workers_per_gather` (integer)

Sets the maximum number of workers that can be started by a single `Gather` or `Gather Merge` node. Parallel workers are taken from the pool of processes established by [max_worker_processes](#), limited by [max_parallel_workers](#). Note that the requested number of workers may not actually be available at run time. If this occurs, the plan will run with fewer workers than expected, which may be inefficient. The default value is 2. Setting this value to 0 disables parallel query execution.

Note that parallel queries may consume very substantially more resources than non-parallel queries, because each worker process is a completely separate process which has roughly the same impact on the system as an additional user session. This should be taken into account when choosing a value for this setting, as well as when configuring other settings that control resource utilization, such as [work_mem](#). Resource limits such as `work_mem` are applied individually to each worker, which means the total utilization may be much higher across all processes than it would normally be for any single process. For example, a parallel query using 4 workers may use up to 5 times as much CPU time, memory, I/O bandwidth, and so forth as a query which uses no workers at all.

For more information on parallel query, see [Chapter 15](#).

`max_parallel_maintenance_workers` (integer)

Sets the maximum number of parallel workers that can be started by a single utility command. Currently, the parallel utility commands that support the use of parallel workers are `CREATE INDEX` when building a B-tree or BRIN index, and `VACUUM` without `FULL` option. Parallel workers are taken from the pool of processes established by [max_worker_processes](#), limited by [max_parallel_workers](#). Note that the requested number of workers may not actually be available at run time. If this occurs, the utility operation will run with fewer workers than expected. The default value is 2. Setting this value to 0 disables the use of parallel workers by utility commands.

Note that parallel utility commands should not consume substantially more memory than equivalent non-parallel operations. This strategy differs from that of parallel query, where resource limits generally apply per worker process. Parallel utility commands treat the resource limit `main-`

`nance_work_mem` as a limit to be applied to the entire utility command, regardless of the number of parallel worker processes. However, parallel utility commands may still consume substantially more CPU resources and I/O bandwidth.

`max_parallel_workers (integer)`

Sets the maximum number of workers that the cluster can support for parallel operations. The default value is 8. When increasing or decreasing this value, consider also adjusting [max_parallel_maintenance_workers](#) and [max_parallel_workers_per_gather](#). Also, note that a setting for this value which is higher than [max_worker_processes](#) will have no effect, since parallel workers are taken from the pool of worker processes established by that setting.

`parallel_leader_participation (boolean)`

Allows the leader process to execute the query plan under `Gather` and `Gather Merge` nodes instead of waiting for worker processes. The default is `on`. Setting this value to `off` reduces the likelihood that workers will become blocked because the leader is not reading tuples fast enough, but requires the leader process to wait for worker processes to start up before the first tuples can be produced. The degree to which the leader can help or hinder performance depends on the plan type, number of workers and query duration.

19.4.7. Prioritization

`usage_tracking_interval (integer)`

Sets the time interval, in seconds, for calculating usage statistics. Based on this statistics, Postgres Pro Shardman can control resource usage for each session in accordance with the prioritization policy configured by [session_cpu_weight](#), [session_ioread_weight](#), and [session_iowrite_weight](#) parameters.

When set to a positive value, this parameter starts a background worker that collects statistics on CPU time, the number of local and shared blocks read by the backends, and the number of local and shared blocks dirtied by the backends. Avoid setting this parameter to a small value as frequent statistic collection can cause overhead.

The default value is zero, which disables statistics collection and resource prioritization.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`session_cpu_weight (integer)`

Sets CPU usage weight for the current session. Possible values are 1, 2, 4, and 8. The higher the value, the more resources the session can use as compared to sessions with lower weights.

The [usage_tracking_interval](#) parameter must be set to a positive value for this setting to take effect. Resource usage is planned based on the statistics collected for the previous interval defined by [usage_tracking_interval](#). If all sessions have the same weight, Postgres Pro Shardman does not prioritize resource usage.

Default: 4

`session_ioread_weight (integer)`

Sets the weight for reading local and shared blocks for the current session. Possible values are 1, 2, 4, and 8. The higher the value, the more resources the session can use as compared to sessions with lower weights.

The [usage_tracking_interval](#) parameter must be set to a positive value for this setting to take effect. Resource usage is planned based on the statistics collected for the previous interval defined by [usage_tracking_interval](#). If all sessions have the same weight, Postgres Pro Shardman does not prioritize resource usage.

Default: 4

`session_iowrite_weight` (integer)

Sets the weight for writing to local and shared blocks for the current session. Possible values are 1, 2, 4, and 8. The higher the value, the more resources the session can use as compared to sessions with lower weights.

The [usage_tracking_interval](#) parameter must be set to a positive value for this setting to take effect. Resource usage is planned based on the statistics collected for the previous interval defined by [usage_tracking_interval](#). If all sessions have the same weight, Postgres Pro Shardman does not prioritize resource usage.

Default: 4

19.5. Write Ahead Log

For additional information on tuning these settings, see [Section 28.5](#).

19.5.1. Settings

`wal_level` (enum)

`wal_level` determines how much information is written to the WAL. The default value is `replica`, which writes enough data to support WAL archiving and replication, including running read-only queries on a standby server. `minimal` removes all logging except the information required to recover from a crash or immediate shutdown. Finally, `logical` adds information necessary to support logical decoding. Each level includes the information logged at all lower levels. This parameter can only be set at server start.

For distributed system specifics, see [wal_level](#).

The `minimal` level generates the least WAL volume. It logs no row information for permanent relations in transactions that create or rewrite them. This can make operations much faster. Operations that initiate this optimization include:

```
ALTER ... SET TABLESPACE
CLUSTER
CREATE TABLE
REFRESH MATERIALIZED VIEW (without CONCURRENTLY)
REINDEX
TRUNCATE
```

However, `minimal` WAL does not contain sufficient information for point-in-time recovery, so `replica` or higher must be used to enable continuous archiving ([archive_mode](#)) and streaming binary replication. In fact, the server will not even start in this mode if `max_wal_senders` is non-zero. Note that changing `wal_level` to `minimal` makes previous base backups unusable for point-in-time recovery and standby servers.

In `logical` level, the same information is logged as with `replica`, plus information needed to extract logical change sets from the WAL. Using a level of `logical` will increase the WAL volume, particularly if many tables are configured for `REPLICA IDENTITY FULL` and many `UPDATE` and `DELETE` statements are executed.

In releases prior to 9.6, this parameter also allowed the values `archive` and `hot_standby`. These are still accepted but mapped to `replica`.

`fsync` (boolean)

If this parameter is on, the Postgres Pro server will try to make sure that updates are physically written to disk, by issuing `fsync()` system calls or various equivalent methods (see [wal_sync_method](#)). This ensures that the database cluster can recover to a consistent state after an operating system or hardware crash.

While turning off `fsync` is often a performance benefit, this can result in unrecoverable data corruption in the event of a power failure or system crash. Thus it is only advisable to turn off `fsync` if you can easily recreate your entire database from external data.

Examples of safe circumstances for turning off `fsync` include the initial loading of a new database cluster from a backup file, using a database cluster for processing a batch of data after which the database will be thrown away and recreated, or for a read-only database clone which gets recreated frequently and is not used for failover. High quality hardware alone is not a sufficient justification for turning off `fsync`.

For reliable recovery when changing `fsync` off to on, it is necessary to force all modified buffers in the kernel to durable storage. This can be done while the cluster is shutdown or while `fsync` is on by running `initdb --sync-only`, running `sync`, unmounting the file system, or rebooting the server.

In many situations, turning off `synchronous_commit` for noncritical transactions can provide much of the potential performance benefit of turning off `fsync`, without the attendant risks of data corruption.

`fsync` can only be set in the `postgresql.conf` file or on the server command line. If you turn this parameter off, also consider turning off `full_page_writes`.

`synchronous_commit` (enum)

Specifies how much WAL processing must complete before the database server returns a “success” indication to the client. Valid values are `remote_apply`, `on` (the default), `remote_write`, `local`, and `off`.

If `synchronous_standby_names` is empty, the only meaningful settings are `on` and `off`; `remote_apply`, `remote_write` and `local` all provide the same local synchronization level as `on`. The local behavior of all non-`off` modes is to wait for local flush of WAL to disk. In `off` mode, there is no waiting, so there can be a delay between when success is reported to the client and when the transaction is later guaranteed to be safe against a server crash. (The maximum delay is three times `wal_writer_delay`.) Unlike `fsync`, setting this parameter to `off` does not create any risk of database inconsistency: an operating system or database crash might result in some recent allegedly-committed transactions being lost, but the database state will be just the same as if those transactions had been aborted cleanly. So, turning `synchronous_commit` off can be a useful alternative when performance is more important than exact certainty about the durability of a transaction. For more discussion see [Section 28.4](#).

If `synchronous_standby_names` is non-empty, `synchronous_commit` also controls whether transaction commits will wait for their WAL records to be processed on the standby server(s).

When set to `remote_apply`, commits will wait until replies from the current synchronous standby(s) indicate they have received the commit record of the transaction and applied it, so that it has become visible to queries on the standby(s), and also written to durable storage on the standbys. This will cause much larger commit delays than previous settings since it waits for WAL replay. When set to `on`, commits wait until replies from the current synchronous standby(s) indicate they have received the commit record of the transaction and flushed it to durable storage. This ensures the transaction will not be lost unless both the primary and all synchronous standbys suffer corruption of their database storage. When set to `remote_write`, commits will wait until replies from the current synchronous standby(s) indicate they have received the commit record of the transaction and written it to their file systems. This setting ensures data preservation if a standby instance of Postgres Pro crashes, but not if the standby suffers an operating-system-level crash because the data has not necessarily reached durable storage on the standby. The setting `local` causes commits to wait for local flush to disk, but not for replication. This is usually not desirable when synchronous replication is in use, but is provided for completeness.

This parameter can be changed at any time; the behavior for any one transaction is determined by the setting in effect when it commits. It is therefore possible, and useful, to have some transactions commit synchronously and others asynchronously. For example, to make a single multistatement transaction commit asynchronously when the default is the opposite, issue `SET LOCAL synchronous_commit TO OFF` within the transaction.

Table 19.1 summarizes the capabilities of the `synchronous_commit` settings.

Table 19.1. `synchronous_commit` Modes

synchronous_commit setting	local durable commit	standby durable commit after PG crash	standby durable commit after OS crash	standby query consistency
<code>remote_apply</code>	•	•	•	•
<code>on</code>	•	•	•	
<code>remote_write</code>	•	•		
<code>local</code>	•			
<code>off</code>				

`wal_sync_method` (enum)

Method used for forcing WAL updates out to disk. If `fsync` is off then this setting is irrelevant, since WAL file updates will not be forced out at all. Possible values are:

- `open_datasync` (write WAL files with `open()` option `O_DSYNC`)
- `fdasync` (call `fdasync()` at each commit)
- `fsync` (call `fsync()` at each commit)
- `fsync_writethrough` (call `fsync()` at each commit, forcing write-through of any disk write cache)
- `open_sync` (write WAL files with `open()` option `O_SYNC`)

Not all of these choices are available on all platforms. The default is the first method in the above list that is supported by the platform, except that `fdasync` is the default on Linux and FreeBSD. The default is not necessarily ideal; it might be necessary to change this setting or other aspects of your system configuration in order to create a crash-safe configuration or achieve optimal performance. These aspects are discussed in [Section 28.1](#). This parameter can only be set in the `postgresql.conf` file or on the server command line.

`full_page_writes` (boolean)

When this parameter is on, the Postgres Pro server writes the entire content of each disk page to WAL during the first modification of that page after a checkpoint. This is needed because a page write that is in process during an operating system crash might be only partially completed, leading to an on-disk page that contains a mix of old and new data. The row-level change data normally stored in WAL will not be enough to completely restore such a page during post-crash recovery. Storing the full page image guarantees that the page can be correctly restored, but at the price of increasing the amount of data that must be written to WAL. (Because WAL replay always starts from a checkpoint, it is sufficient to do this during the first change of each page after a checkpoint. Therefore, one way to reduce the cost of full-page writes is to increase the checkpoint interval parameters.)

Turning this parameter off speeds normal operation, but might lead to either unrecoverable data corruption, or silent data corruption, after a system failure. The risks are similar to turning off `fsync`, though smaller, and it should be turned off only based on the same circumstances recommended for that parameter.

Turning off this parameter does not affect use of WAL archiving for point-in-time recovery (PITR) (see [Section 25.1](#)).

This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `on`.

`wal_log_hints (boolean)`

When this parameter is `on`, the Postgres Pro server writes the entire content of each disk page to WAL during the first modification of that page after a checkpoint, even for non-critical modifications of so-called hint bits.

If data checksums are enabled, hint bit updates are always WAL-logged and this setting is ignored. You can use this setting to test how much extra WAL-logging would occur if your database had data checksums enabled.

This parameter can only be set at server start. The default value is `off`.

`wal_compression (enum)`

This parameter enables compression of WAL using the specified compression method. When enabled, the Postgres Pro server compresses full page images written to WAL when [full_page_writes](#) is `on` or during a base backup. A compressed page image will be decompressed during WAL replay. The supported methods are `pglz`, `lz4` (if Postgres Pro was compiled with `--with-lz4`) and `zstd` (if Postgres Pro was compiled with `--with-zstd`). The default value is `lz4`. To disable compression, it can be set to `off` instead. Only superusers and users with the appropriate `SET` privilege can change this setting.

Enabling compression can reduce the WAL volume without increasing the risk of unrecoverable data corruption, but at the cost of some extra CPU spent on the compression during WAL logging and on the decompression during WAL replay.

`wal_init_zero (boolean)`

If set to `on` (the default), this option causes new WAL files to be filled with zeroes. On some file systems, this ensures that space is allocated before we need to write WAL records. However, *Copy-On-Write* (COW) file systems may not benefit from this technique, so the option is given to skip the unnecessary work. If set to `off`, only the final byte is written when the file is created so that it has the expected size.

`wal_recycle (boolean)`

If set to `on` (the default), this option causes WAL files to be recycled by renaming them, avoiding the need to create new ones. On COW file systems, it may be faster to create new ones, so the option is given to disable this behavior.

`wal_buffers (integer)`

The amount of shared memory used for WAL data that has not yet been written to disk. The default setting of `-1` selects a size equal to 1/32nd (about 3%) of [shared_buffers](#), but not less than 64kB nor more than the size of one WAL segment, typically 16MB. This value can be set manually if the automatic choice is too large or too small, but any positive value less than 32kB will be treated as 32kB. If this value is specified without units, it is taken as WAL blocks, that is `XLOG_BLCKSZ` bytes, typically 8kB. This parameter can only be set at server start.

The contents of the WAL buffers are written out to disk at every transaction commit, so extremely large values are unlikely to provide a significant benefit. However, setting this value to at least a few megabytes can improve write performance on a busy server where many clients are committing at once. The auto-tuning selected by the default setting of `-1` should give reasonable results in most cases.

`wal_writer_delay (integer)`

Specifies how often the WAL writer flushes WAL, in time terms. After flushing WAL the writer sleeps for the length of time given by `wal_writer_delay`, unless woken up sooner by an asynchronously committing transaction. If the last flush happened less than `wal_writer_delay` ago and less than `wal_writer_flush_after` worth of WAL has been produced since, then WAL is only written to the operating system, not flushed to disk. If this value is specified without units, it is taken as milliseconds. The default value is 200 milliseconds (200ms). Note that on some systems, the effective reso-

lution of sleep delays is 10 milliseconds; setting `wal_writer_delay` to a value that is not a multiple of 10 might have the same results as setting it to the next higher multiple of 10. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`wal_writer_flush_after` (integer)

Specifies how often the WAL writer flushes WAL, in volume terms. If the last flush happened less than `wal_writer_delay` ago and less than `wal_writer_flush_after` worth of WAL has been produced since, then WAL is only written to the operating system, not flushed to disk. If `wal_writer_flush_after` is set to 0 then WAL data is always flushed immediately. If this value is specified without units, it is taken as WAL blocks, that is `XLOG_BLCKSZ` bytes, typically 8kB. The default is 1MB. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`wal_skip_threshold` (integer)

When `wal_level` is `minimal` and a transaction commits after creating or rewriting a permanent relation, this setting determines how to persist the new data. If the data is smaller than this setting, write it to the WAL log; otherwise, use an `fsync` of affected files. Depending on the properties of your storage, raising or lowering this value might help if such commits are slowing concurrent transactions. If this value is specified without units, it is taken as kilobytes. The default is two megabytes (2MB).

`commit_delay` (integer)

Setting `commit_delay` adds a time delay before a WAL flush is initiated. This can improve group commit throughput by allowing a larger number of transactions to commit via a single WAL flush, if system load is high enough that additional transactions become ready to commit within the given interval. However, it also increases latency by up to the `commit_delay` for each WAL flush. Because the delay is just wasted if no other transactions become ready to commit, a delay is only performed if at least `commit_siblings` other transactions are active when a flush is about to be initiated. Also, no delays are performed if `fsync` is disabled. If this value is specified without units, it is taken as microseconds. The default `commit_delay` is zero (no delay). Only superusers and users with the appropriate `SET` privilege can change this setting.

In PostgreSQL releases prior to 9.3, `commit_delay` behaved differently and was much less effective: it affected only commits, rather than all WAL flushes, and waited for the entire configured delay even if the WAL flush was completed sooner. Beginning in PostgreSQL 9.3, the first process that becomes ready to flush waits for the configured interval, while subsequent processes wait only until the leader completes the flush operation.

`commit_siblings` (integer)

Minimum number of concurrent open transactions to require before performing the `commit_delay` delay. A larger value makes it more probable that at least one other transaction will become ready to commit during the delay interval. The default is five transactions.

`wal_sender_check_crc` (boolean)

When this parameter is `on`, during replication the WAL sender process checks the CRC value of each record before sending it to a replica, and memory for a second copy of [WAL buffers](#) is allocated. The default value is `off`. This parameter can only be set at server start.

`wal_sender_panic_on_crc_error` (boolean)

When this parameter is `on`, Postgres Pro raises a PANIC-level error if a corrupted WAL record could not be restored from WAL buffers. The default value is `off` for raising a FATAL-level error. This parameter can only be set at server start.

19.5.2. Checkpoints

`checkpoint_timeout` (integer)

Maximum time between automatic WAL checkpoints. If this value is specified without units, it is taken as seconds. The valid range is between 30 seconds and one day. The default is five minutes

(5min). Increasing this parameter can increase the amount of time needed for crash recovery. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`checkpoint_completion_target` (floating point)

Specifies the target of checkpoint completion, as a fraction of total time between checkpoints. The default is 0.9, which spreads the checkpoint across almost all of the available interval, providing fairly consistent I/O load while also leaving some time for checkpoint completion overhead. Reducing this parameter is not recommended because it causes the checkpoint to complete faster. This results in a higher rate of I/O during the checkpoint followed by a period of less I/O between the checkpoint completion and the next scheduled checkpoint. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`checkpoint_flush_after` (integer)

Whenever more than this amount of data has been written while performing a checkpoint, attempt to force the OS to issue these writes to the underlying storage. Doing so will limit the amount of dirty data in the kernel's page cache, reducing the likelihood of stalls when an `fsync` is issued at the end of the checkpoint, or when the OS writes data back in larger batches in the background. Often that will result in greatly reduced transaction latency, but there also are some cases, especially with workloads that are bigger than [shared_buffers](#), but smaller than the OS's page cache, where performance might degrade. This setting may have no effect on some platforms. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The valid range is between 0, which disables forced writeback, and 2MB. The default is 256kB on Linux, 0 elsewhere. (If `BLCKSZ` is not 8kB, the default and maximum values scale proportionally to it.) This parameter can only be set in the `postgresql.conf` file or on the server command line.

`checkpoint_warning` (integer)

Write a message to the server log if checkpoints caused by the filling of WAL segment files happen closer together than this amount of time (which suggests that `max_wal_size` ought to be raised). If this value is specified without units, it is taken as seconds. The default is 30 seconds (30s). Zero disables the warning. No warnings will be generated if `checkpoint_timeout` is less than `checkpoint_warning`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`max_wal_size` (integer)

Maximum size to let the WAL grow during automatic checkpoints. This is a soft limit; WAL size can exceed `max_wal_size` under special circumstances, such as heavy load, a failing `archive_command` or `archive_library`, or a high `wal_keep_size` setting. If this value is specified without units, it is taken as megabytes. The default is 1 GB. Increasing this parameter can increase the amount of time needed for crash recovery. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`min_wal_size` (integer)

As long as WAL disk usage stays below this setting, old WAL files are always recycled for future use at a checkpoint, rather than removed. This can be used to ensure that enough WAL space is reserved to handle spikes in WAL usage, for example when running large batch jobs. If this value is specified without units, it is taken as megabytes. The default is 80 MB. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.5.3. Archiving

`archive_mode` (enum)

When `archive_mode` is enabled, completed WAL segments are sent to archive storage by setting [archive_command](#) or [archive_library](#). In addition to `off`, to disable, there are two modes: `on`, and `always`. During normal operation, there is no difference between the two modes, but when set to `always` the WAL archiver is enabled also during archive recovery or standby mode. In `always` mode,

all files restored from the archive or streamed with streaming replication will be archived (again). See [Section 26.2.9](#) for details.

`archive_mode` is a separate setting from `archive_command` and `archive_library` so that `archive_command` and `archive_library` can be changed without leaving archiving mode. This parameter can only be set at server start. `archive_mode` cannot be enabled when `wal_level` is set to `minimal`.

`archive_command` (string)

The local shell command to execute to archive a completed WAL file segment. Any `%p` in the string is replaced by the path name of the file to archive, and any `%f` is replaced by only the file name. (The path name is relative to the working directory of the server, i.e., the cluster's data directory.) Use `%` to embed an actual `%` character in the command. It is important for the command to return a zero exit status only if it succeeds. For more information see [Section 25.1.1](#).

This parameter can only be set in the `postgresql.conf` file or on the server command line. It is only used if `archive_mode` was enabled at server start and `archive_library` is set to an empty string. If both `archive_command` and `archive_library` are set, an error will be raised. If `archive_command` is an empty string (the default) while `archive_mode` is enabled (and `archive_library` is set to an empty string), WAL archiving is temporarily disabled, but the server continues to accumulate WAL segment files in the expectation that a command will soon be provided. Setting `archive_command` to a command that does nothing but return true, e.g., `/bin/true` (REM on Windows), effectively disables archiving, but also breaks the chain of WAL files needed for archive recovery, so it should only be used in unusual circumstances.

`archive_library` (string)

The library to use for archiving completed WAL file segments. If set to an empty string (the default), archiving via shell is enabled, and [archive_command](#) is used. If both `archive_command` and `archive_library` are set, an error will be raised. Otherwise, the specified shared library is used for archiving. The WAL archiver process is restarted by the postmaster when this parameter changes. For more information, see [Section 25.1.1](#) and [Chapter 52](#).

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`archive_timeout` (integer)

The [archive_command](#) or [archive_library](#) is only invoked for completed WAL segments. Hence, if your server generates little WAL traffic (or has slack periods where it does so), there could be a long delay between the completion of a transaction and its safe recording in archive storage. To limit how old unarchived data can be, you can set `archive_timeout` to force the server to switch to a new WAL segment file periodically. When this parameter is greater than zero, the server will switch to a new segment file whenever this amount of time has elapsed since the last segment file switch, and there has been any database activity, including a single checkpoint (checkpoints are skipped if there is no database activity). Note that archived files that are closed early due to a forced switch are still the same length as completely full files. Therefore, it is unwise to use a very short `archive_timeout` — it will bloat your archive storage. `archive_timeout` settings of a minute or so are usually reasonable. You should consider using streaming replication, instead of archiving, if you want data to be copied off the primary server more quickly than that. If this value is specified without units, it is taken as seconds. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.5.4. Recovery

This section describes the settings that apply to recovery in general, affecting crash recovery, streaming replication and archive-based replication.

`recovery_prefetch` (enum)

Whether to try to prefetch blocks that are referenced in the WAL that are not yet in the buffer pool, during recovery. Valid values are `off`, `on` and `try` (the default). The setting `try` enables prefetching

only if the operating system provides the `posix_fadvise` function, which is currently used to implement prefetching. Note that some operating systems provide the function, but it doesn't do anything.

Prefetching blocks that will soon be needed can reduce I/O wait times during recovery with some workloads. See also the [wal_decode_buffer_size](#) and [maintenance_io_concurrency](#) settings, which limit prefetching activity.

`wal_decode_buffer_size` (integer)

A limit on how far ahead the server can look in the WAL, to find blocks to prefetch. If this value is specified without units, it is taken as bytes. The default is 512kB.

19.5.5. Archive Recovery

This section describes the settings that apply only for the duration of the recovery. They must be reset for any subsequent recovery you wish to perform.

“Recovery” covers using the server as a standby or for executing a targeted recovery. Typically, standby mode would be used to provide high availability and/or read scalability, whereas a targeted recovery is used to recover from data loss.

To start the server in standby mode, create a file called `standby.signal` in the data directory. The server will enter recovery and will not stop recovery when the end of archived WAL is reached, but will keep trying to continue recovery by connecting to the sending server as specified by the `primary_conninfo` setting and/or by fetching new WAL segments using `restore_command`. For this mode, the parameters from this section and [Section 19.6.3](#) are of interest. Parameters from [Section 19.5.6](#) will also be applied but are typically not useful in this mode.

To start the server in targeted recovery mode, create a file called `recovery.signal` in the data directory. If both `standby.signal` and `recovery.signal` files are created, standby mode takes precedence. Targeted recovery mode ends when the archived WAL is fully replayed, or when `recovery_target` is reached. In this mode, the parameters from both this section and [Section 19.5.6](#) will be used.

`restore_command` (string)

The local shell command to execute to retrieve an archived segment of the WAL file series. This parameter is required for archive recovery, but optional for streaming replication. Any `%f` in the string is replaced by the name of the file to retrieve from the archive, and any `%p` is replaced by the copy destination path name on the server. (The path name is relative to the current working directory, i.e., the cluster's data directory.) Any `%r` is replaced by the name of the file containing the last valid restart point. That is the earliest file that must be kept to allow a restore to be restartable, so this information can be used to truncate the archive to just the minimum required to support restarting from the current restore. `%r` is typically only used by warm-standby configurations (see [Section 26.2](#)). Write `%%` to embed an actual `%` character.

It is important for the command to return a zero exit status only if it succeeds. The command *will* be asked for file names that are not present in the archive; it must return nonzero when so asked. Examples:

```
restore_command = 'cp /mnt/server/archivedir/%f "%p"'
restore_command = 'copy "C:\\server\\archivedir\\%f" "%p"' # Windows
```

An exception is that if the command was terminated by a signal (other than SIGTERM, which is used as part of a database server shutdown) or an error by the shell (such as command not found), then recovery will abort and the server will not start up.

`archive_cleanup_command` (string)

This optional parameter specifies a shell command that will be executed at every restartpoint. The purpose of `archive_cleanup_command` is to provide a mechanism for cleaning up old archived WAL files that are no longer needed by the standby server. Any `%r` is replaced by the name of the file containing the last valid restart point. That is the earliest file that must be *kept* to allow a restore to

be restartable, and so all files earlier than %r may be safely removed. This information can be used to truncate the archive to just the minimum required to support restart from the current restore. The [pg_archivecleanup](#) module is often used in `archive_cleanup_command` for single-standby configurations, for example:

```
archive_cleanup_command = 'pg_archivecleanup /mnt/server/archivedir %r'
```

Note however that if multiple standby servers are restoring from the same archive directory, you will need to ensure that you do not delete WAL files until they are no longer needed by any of the servers. `archive_cleanup_command` would typically be used in a warm-standby configuration (see [Section 26.2](#)). Write %% to embed an actual % character in the command.

If the command returns a nonzero exit status then a warning log message will be written. An exception is that if the command was terminated by a signal or an error by the shell (such as command not found), a fatal error will be raised.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`recovery_end_command (string)`

This parameter specifies a shell command that will be executed once only at the end of recovery. This parameter is optional. The purpose of the `recovery_end_command` is to provide a mechanism for cleanup following replication or recovery. Any %r is replaced by the name of the file containing the last valid restart point, like in [archive_cleanup_command](#).

If the command returns a nonzero exit status then a warning log message will be written and the database will proceed to start up anyway. An exception is that if the command was terminated by a signal or an error by the shell (such as command not found), the database will not proceed with startup.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.5.6. Recovery Target

By default, recovery will recover to the end of the WAL log. The following parameters can be used to specify an earlier stopping point. At most one of `recovery_target`, `recovery_target_lsn`, `recovery_target_name`, `recovery_target_time`, or `recovery_target_xid` can be used; if more than one of these is specified in the configuration file, an error will be raised. These parameters can only be set at server start.

`recovery_target = 'immediate'`

This parameter specifies that recovery should end as soon as a consistent state is reached, i.e., as early as possible. When restoring from an online backup, this means the point where taking the backup ended.

Technically, this is a string parameter, but 'immediate' is currently the only allowed value.

`recovery_target_name (string)`

This parameter specifies the named restore point (created with `pg_create_restore_point()`) to which recovery will proceed.

`recovery_target_time (timestamp)`

This parameter specifies the time stamp up to which recovery will proceed. The precise stopping point is also influenced by [recovery_target_inclusive](#).

The value of this parameter is a time stamp in the same format accepted by the `timestamp with time zone` data type, except that you cannot use a time zone abbreviation (unless the [timezone_abbreviations](#) variable has been set earlier in the configuration file). Preferred style is to use a numeric offset from UTC, or you can write a full time zone name, e.g., `Europe/Helsinki` not `EEST`.

`recovery_target_xid` (string)

This parameter specifies the transaction ID up to which recovery will proceed. Keep in mind that while transaction IDs are assigned sequentially at transaction start, transactions can complete in a different numeric order. The transactions that will be recovered are those that committed before (and optionally including) the specified one. The precise stopping point is also influenced by [recovery_target_inclusive](#).

`recovery_target_lsn` (pg_lsn)

This parameter specifies the LSN of the write-ahead log location up to which recovery will proceed. The precise stopping point is also influenced by [recovery_target_inclusive](#). This parameter is parsed using the system data type `pg_lsn`.

The following options further specify the recovery target, and affect what happens when the target is reached:

`recovery_target_inclusive` (boolean)

Specifies whether to stop just after the specified recovery target (`on`), or just before the recovery target (`off`). Applies when [recovery_target_lsn](#), [recovery_target_time](#), or [recovery_target_xid](#) is specified. This setting controls whether transactions having exactly the target WAL location (LSN), commit time, or transaction ID, respectively, will be included in the recovery. Default is `on`.

`recovery_target_timeline` (string)

Specifies recovering into a particular timeline. The value can be a numeric timeline ID or a special value. The value `current` recovers along the same timeline that was current when the base backup was taken. The value `latest` recovers to the latest timeline found in the archive, which is useful in a standby server. `latest` is the default.

To specify a timeline ID in hexadecimal (for example, if extracted from a WAL file name or history file), prefix it with a `0x`. For instance, if the WAL file name is `00000011000000A10000004F`, then the timeline ID is `0x11` (or 17 decimal).

You usually only need to set this parameter in complex re-recovery situations, where you need to return to a state that itself was reached after a point-in-time recovery. See [Section 25.1.2](#) for discussion.

`recovery_target_action` (enum)

Specifies what action the server should take once the recovery target is reached. The default is `pause`, which means recovery will be paused. `promote` means the recovery process will finish and the server will start to accept connections. Finally `shutdown` will stop the server after reaching the recovery target.

The intended use of the `pause` setting is to allow queries to be executed against the database to check if this recovery target is the most desirable point for recovery. The paused state can be resumed by using `pg_wal_replay_resume()` (see [Table 9.97](#)), which then causes recovery to end. If this recovery target is not the desired stopping point, then shut down the server, change the recovery target settings to a later target and restart to continue recovery.

The `shutdown` setting is useful to have the instance ready at the exact replay point desired. The instance will still be able to replay more WAL records (and in fact will have to replay WAL records since the last checkpoint next time it is started).

Note that because `recovery.signal` will not be removed when `recovery_target_action` is set to `shutdown`, any subsequent start will end with immediate shutdown unless the configuration is changed or the `recovery.signal` file is removed manually.

This setting has no effect if no recovery target is set. If [hot_standby](#) is not enabled, a setting of `pause` will act the same as `shutdown`. If the recovery target is reached while a promotion is ongoing, a setting of `pause` will act the same as `promote`.

In any case, if a recovery target is configured but the archive recovery ends before the target is reached, the server will shut down with a fatal error.

19.5.7. WAL Summarization

These settings control WAL summarization, a feature which must be enabled in order to perform an incremental backup.

`summarize_wal` (boolean)

Enables the WAL summarizer process. Note that WAL summarization can be enabled either on a primary or on a standby. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `off`.

The server cannot be started with `summarize_wal=on` if `wal_level` is set to `minimal`. If `summarize_wal=on` is configured after server startup while `wal_level=minimal`, the summarizer will run but refuse to generate summary files for any WAL generated with `wal_level=minimal`.

`wal_summary_keep_time` (integer)

Configures the amount of time after which the WAL summarizer automatically removes old WAL summaries. The file timestamp is used to determine which files are old enough to remove. Typically, you should set this comfortably higher than the time that could pass between a backup and a later incremental backup that depends on it. WAL summaries must be available for the entire range of WAL records between the preceding backup and the new one being taken; if not, the incremental backup will fail. If this parameter is set to zero, WAL summaries will not be automatically deleted, but it is safe to manually remove files that you know will not be required for future incremental backups. This parameter can only be set in the `postgresql.conf` file or on the server command line. If this value is specified without units, it is taken as minutes. The default is 10 days. If `summarize_wal = off`, existing WAL summaries will not be removed regardless of the value of this parameter, because the WAL summarizer will not run.

19.6. Replication

These settings control the behavior of the built-in *streaming replication* feature (see [Section 26.2.5](#)), and the built-in *logical replication* feature (see [Chapter 29](#)).

For *streaming replication*, servers will be either a primary or a standby server. Primaries can send data, while standbys are always receivers of replicated data. When cascading replication (see [Section 26.2.7](#)) is used, standby servers can also be senders, as well as receivers. Parameters are mainly for sending and standby servers, though some parameters have meaning only on the primary server. Settings may vary across the cluster without problems if that is required.

For *logical replication*, *publishers* (servers that do `CREATE PUBLICATION`) replicate data to *subscribers* (servers that do `CREATE SUBSCRIPTION`). Servers can also be publishers and subscribers at the same time. Note, the following sections refer to publishers as "senders". For more details about logical replication configuration settings refer to [Section 29.11](#).

19.6.1. Sending Servers

These parameters can be set on any server that is to send replication data to one or more standby servers. The primary is always a sending server, so these parameters must always be set on the primary. The role and meaning of these parameters does not change after a standby becomes the primary.

`max_wal_senders` (integer)

Specifies the maximum number of concurrent connections from standby servers or streaming base backup clients (i.e., the maximum number of simultaneously running WAL sender processes). The default is 10. The value 0 means replication is disabled. Abrupt disconnection of a streaming client might leave an orphaned connection slot behind until a timeout is reached, so this parameter should be set slightly higher than the maximum number of expected clients so disconnected clients can

immediately reconnect. This parameter can only be set at server start. Also, `wal_level` must be set to `replica` or higher to allow connections from standby servers.

When running a standby server, you must set this parameter to the same or higher value than on the primary server. Otherwise, queries will not be allowed in the standby server.

`max_replication_slots` (integer)

Specifies the maximum number of replication slots (see [Section 26.2.6](#)) that the server can support. The default is 10. This parameter can only be set at server start. Setting it to a lower value than the number of currently existing replication slots will prevent the server from starting. Also, `wal_level` must be set to `replica` or higher to allow replication slots to be used.

Note that this parameter also applies on the subscriber side, but with a different meaning.

`wal_keep_size` (integer)

Specifies the minimum size of past WAL files kept in the `pg_wal` directory, in case a standby server needs to fetch them for streaming replication. If a standby server connected to the sending server falls behind by more than `wal_keep_size` megabytes, the sending server might remove a WAL segment still needed by the standby, in which case the replication connection will be terminated. Downstream connections will also eventually fail as a result. (However, the standby server can recover by fetching the segment from archive, if WAL archiving is in use.)

This sets only the minimum size of segments retained in `pg_wal`; the system might need to retain more segments for WAL archival or to recover from a checkpoint. If `wal_keep_size` is zero (the default), the system doesn't keep any extra segments for standby purposes, so the number of old WAL segments available to standby servers is a function of the location of the previous checkpoint and status of WAL archiving. If this value is specified without units, it is taken as megabytes. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`max_slot_wal_keep_size` (integer)

Specify the maximum size of WAL files that [replication slots](#) are allowed to retain in the `pg_wal` directory at checkpoint time. If `max_slot_wal_keep_size` is -1 (the default), replication slots may retain an unlimited amount of WAL files. Otherwise, if `restart_lsn` of a replication slot falls behind the current LSN by more than the given size, the standby using the slot may no longer be able to continue replication due to removal of required WAL files. You can see the WAL availability of replication slots in [pg_replication_slots](#). If this value is specified without units, it is taken as megabytes. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`wal_sender_timeout` (integer)

Terminate replication connections that are inactive for longer than this amount of time. This is useful for the sending server to detect a standby crash or network outage. If this value is specified without units, it is taken as milliseconds. The default value is 60 seconds. A value of zero disables the timeout mechanism.

With a cluster distributed across multiple geographic locations, using different values per location brings more flexibility in the cluster management. A smaller value is useful for faster failure detection with a standby having a low-latency network connection, and a larger value helps in judging better the health of a standby if located on a remote location, with a high-latency network connection.

`track_commit_timestamp` (boolean)

Record commit time of transactions. This parameter can only be set in `postgresql.conf` file or on the server command line. The default value is `off`.

`page_repair` (boolean)

Enable page repair via streaming replication from standby in case of data corruption. This parameter can only be set in `postgresql.conf` file or on the server command line. The default value is `off`.

`synchronized_standby_slots` (string)

A comma-separated list of streaming replication standby server slot names that logical WAL sender processes will wait for. Logical WAL sender processes will send decoded changes to plugins only after the specified replication slots confirm receiving WAL. This guarantees that logical replication failover slots do not consume changes until those changes are received and flushed to corresponding physical standbys. If a logical replication connection is meant to switch to a physical standby after the standby is promoted, the physical replication slot for the standby should be listed here. Note that logical replication will not proceed if the slots specified in the `synchronized_standby_slots` do not exist or are invalidated. Additionally, the replication management functions `pg_replication_slot_advance`, `pg_logical_slot_get_changes`, and `pg_logical_slot_peek_changes`, when used with logical failover slots, will block until all physical slots specified in `synchronized_standby_slots` have confirmed WAL receipt.

The standbys corresponding to the physical replication slots in `synchronized_standby_slots` must configure `sync_replication_slots = true` so they can receive logical failover slot changes from the primary.

19.6.2. Primary Server

These parameters can be set on the primary server that is to send replication data to one or more standby servers. Note that in addition to these parameters, `wal_level` must be set appropriately on the primary server, and optionally WAL archiving can be enabled as well (see [Section 19.5.3](#)). The values of these parameters on standby servers are irrelevant, although you may wish to set them there in preparation for the possibility of a standby becoming the primary.

`synchronous_standby_names` (string)

Specifies a list of standby servers that can support *synchronous replication*, as described in [Section 26.2.8](#). There will be one or more active synchronous standbys; transactions waiting for commit will be allowed to proceed after these standby servers confirm receipt of their data. The synchronous standbys will be those whose names appear in this list, and that are both currently connected and streaming data in real-time (as shown by a state of `streaming` in the `pg_stat_replication` view). Specifying more than one synchronous standby can allow for very high availability and protection against data loss.

The name of a standby server for this purpose is the `application_name` setting of the standby, as set in the standby's connection information. In case of a physical replication standby, this should be set in the `primary_conninfo` setting; the default is the setting of `cluster_name` if set, else `walreceiver`. For logical replication, this can be set in the connection information of the subscription, and it defaults to the subscription name. For other replication stream consumers, consult their documentation.

This parameter specifies a list of standby servers using either of the following syntaxes:

```
[FIRST] num_sync [MIN num_sync_min] ( standby_name [, ...] )
ANY num_sync [MIN num_sync_min] ( standby_name [, ...] )
standby_name [, ...]
```

where `num_sync` is the number of synchronous standbys that transactions need to wait for replies from, `num_sync_min` is the minimum number of synchronous standbys that must be connected to the primary server for it to continue running, and `standby_name` is the name of a standby server. `FIRST` and `ANY` specify the method to choose synchronous standbys from the listed servers. `MIN` allows to avoid the primary server blocking if some of the standbys are temporarily unavailable.

The keyword `FIRST`, coupled with `num_sync`, specifies a priority-based synchronous replication and makes transaction commits wait until their WAL records are replicated to `num_sync` synchronous standbys chosen based on their priorities. For example, a setting of `FIRST 3 (s1, s2, s3, s4)` will cause each commit to wait for replies from three higher-priority standbys chosen from standby servers `s1`, `s2`, `s3` and `s4`. The standbys whose names appear earlier in the list are given higher priority and will be considered as synchronous. Other standby servers appearing later in this list represent potential synchronous standbys. If any of the current synchronous standbys disconnects for

whatever reason, it will be replaced immediately with the next-highest-priority standby. The keyword `FIRST` is optional.

The keyword `ANY`, coupled with `num_sync`, specifies a quorum-based synchronous replication and makes transaction commits wait until their WAL records are replicated to *at least* `num_sync` listed standbys. For example, a setting of `ANY 3 (s1, s2, s3, s4)` will cause each commit to proceed as soon as at least any three standbys of `s1`, `s2`, `s3` and `s4` reply.

`FIRST` and `ANY` are case-insensitive. If these keywords are used as the name of a standby server, its `standby_name` must be double-quoted.

The optional `MIN num_sync_min` setting relaxes the standby requirements specified by `FIRST` and `ANY`. If the primary loses connection to one or more synchronous standbys, it continues running while at least `num_sync_min` standby servers are available. Once the connection is restored, the behavior of the primary server depends on the [synchronous_standby_gap](#) setting.

The third syntax was used before PostgreSQL version 9.6 and is still supported. It's the same as the first syntax with `FIRST` and `num_sync` equal to 1. For example, `FIRST 1 (s1, s2)` and `s1, s2` have the same meaning: either `s1` or `s2` is chosen as a synchronous standby.

The special entry `*` matches any standby name.

There is no mechanism to enforce uniqueness of standby names. In case of duplicates one of the matching standbys will be considered as higher priority, though exactly which one is indeterminate.

Note

Each `standby_name` should have the form of a valid SQL identifier, unless it is `*`. You can use double-quoting if necessary. But note that `standby_names` are compared to standby application names case-insensitively, whether double-quoted or not.

If no synchronous standby names are specified here, then synchronous replication is not enabled and transaction commits will not wait for replication. This is the default configuration. Even when synchronous replication is enabled, individual transactions can be configured not to wait for replication by setting the [synchronous_commit](#) parameter to `local` or `off`.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`synchronous_standby_gap` (integer)

Specifies the WAL data gap between the primary server and its synchronous standbys, in kilobytes. Setting this parameter to a positive value allows to minimize the primary server downtime when using synchronous replication. If a standby lags behind the primary for more than the specified gap, the primary server continues running while the standby tries to catch up. Once the standby reaches the specified gap, the primary is blocked to complete standby synchronization. See [Section 26.2.8.1](#) for details.

The default value is 0, which means that the primary is always blocked while a synchronous standby replays WAL data.

19.6.3. Standby Servers

These settings control the behavior of a [standby server](#) that is to receive replication data. Their values on the primary server are irrelevant.

`primary_conninfo` (string)

Specifies a connection string to be used for the standby server to connect with a sending server. This string is in the format described in [Section 35.1.1](#). If any option is unspecified in this string, then

the corresponding environment variable (see [Section 35.15](#)) is checked. If the environment variable is not set either, then defaults are used.

The connection string should specify the host name (or address) of the sending server, as well as the port number if it is not the same as the standby server's default. Also specify a user name corresponding to a suitably-privileged role on the sending server (see [Section 26.2.5.1](#)). A password needs to be provided too, if the sender demands password authentication. It can be provided in the `primary_conninfo` string, or in a separate `~/.pgpass` file on the standby server (use `replication` as the database name).

For replication slot synchronization (see [Section 50.2.3](#)), it is also necessary to specify a valid `dbname` in the `primary_conninfo` string. This will only be used for slot synchronization. It is ignored for streaming.

This parameter can only be set in the `postgresql.conf` file or on the server command line. If this parameter is changed while the WAL receiver process is running, that process is signaled to shut down and expected to restart with the new setting (except if `primary_conninfo` is an empty string). This setting has no effect if the server is not in standby mode.

`primary_slot_name` (string)

Optionally specifies an existing replication slot to be used when connecting to the sending server via streaming replication to control resource removal on the upstream node (see [Section 26.2.6](#)). This parameter can only be set in the `postgresql.conf` file or on the server command line. If this parameter is changed while the WAL receiver process is running, that process is signaled to shut down and expected to restart with the new setting. This setting has no effect if `primary_conninfo` is not set or the server is not in standby mode.

`hot_standby` (boolean)

Specifies whether or not you can connect and run queries during recovery, as described in [Section 26.4](#). The default value is `on`. This parameter can only be set at server start. It only has effect during archive recovery or in standby mode.

`max_standby_archive_delay` (integer)

When hot standby is active, this parameter determines how long the standby server should wait before canceling standby queries that conflict with about-to-be-applied WAL entries, as described in [Section 26.4.2](#). `max_standby_archive_delay` applies when WAL data is being read from WAL archive (and is therefore not current). If this value is specified without units, it is taken as milliseconds. The default is 30 seconds. A value of -1 allows the standby to wait forever for conflicting queries to complete. This parameter can only be set in the `postgresql.conf` file or on the server command line.

Note that `max_standby_archive_delay` is not the same as the maximum length of time a query can run before cancellation; rather it is the maximum total time allowed to apply any one WAL segment's data. Thus, if one query has resulted in significant delay earlier in the WAL segment, subsequent conflicting queries will have much less grace time.

`max_standby_streaming_delay` (integer)

When hot standby is active, this parameter determines how long the standby server should wait before canceling standby queries that conflict with about-to-be-applied WAL entries, as described in [Section 26.4.2](#). `max_standby_streaming_delay` applies when WAL data is being received via streaming replication. If this value is specified without units, it is taken as milliseconds. The default is 30 seconds. A value of -1 allows the standby to wait forever for conflicting queries to complete. This parameter can only be set in the `postgresql.conf` file or on the server command line.

Note that `max_standby_streaming_delay` is not the same as the maximum length of time a query can run before cancellation; rather it is the maximum total time allowed to apply WAL data once it has been received from the primary server. Thus, if one query has resulted in significant delay,

subsequent conflicting queries will have much less grace time until the standby server has caught up again.

`wal_receiver_create_temp_slot` (boolean)

Specifies whether the WAL receiver process should create a temporary replication slot on the remote instance when no permanent replication slot to use has been configured (using [primary_slot_name](#)). The default is off. This parameter can only be set in the `postgresql.conf` file or on the server command line. If this parameter is changed while the WAL receiver process is running, that process is signaled to shut down and expected to restart with the new setting.

`wal_receiver_status_interval` (integer)

Specifies the minimum frequency for the WAL receiver process on the standby to send information about replication progress to the primary or upstream standby, where it can be seen using the [pg_stat_replication](#) view. The standby will report the last write-ahead log location it has written, the last position it has flushed to disk, and the last position it has applied. This parameter's value is the maximum amount of time between reports. Updates are sent each time the write or flush positions change, or as often as specified by this parameter if set to a non-zero value. There are additional cases where updates are sent while ignoring this parameter; for example, when processing of the existing WAL completes or when `synchronous_commit` is set to `remote_apply`. Thus, the apply position may lag slightly behind the true position. If this value is specified without units, it is taken as seconds. The default value is 10 seconds. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`hot_standby_feedback` (boolean)

Specifies whether or not a hot standby will send feedback to the primary or upstream standby about queries currently executing on the standby. This parameter can be used to eliminate query cancels caused by cleanup records, but can cause database bloat on the primary for some workloads. Feedback messages will not be sent more frequently than once per `wal_receiver_status_interval`. The default value is `off`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

If cascaded replication is in use the feedback is passed upstream until it eventually reaches the primary. Standbys make no other use of feedback they receive other than to pass upstream.

`wal_receiver_timeout` (integer)

Terminate replication connections that are inactive for longer than this amount of time. This is useful for the receiving standby server to detect a primary node crash or network outage. If this value is specified without units, it is taken as milliseconds. The default value is 60 seconds. A value of zero disables the timeout mechanism. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`wal_retrieve_retry_interval` (integer)

Specifies how long the standby server should wait when WAL data is not available from any sources (streaming replication, local `pg_wal` or WAL archive) before trying again to retrieve WAL data. If this value is specified without units, it is taken as milliseconds. The default value is 5 seconds. This parameter can only be set in the `postgresql.conf` file or on the server command line.

This parameter is useful in configurations where a node in recovery needs to control the amount of time to wait for new WAL data to be available. For example, in archive recovery, it is possible to make the recovery more responsive in the detection of a new WAL file by reducing the value of this parameter. On a system with low WAL activity, increasing it reduces the amount of requests necessary to access WAL archives, something useful for example in cloud environments where the number of times an infrastructure is accessed is taken into account.

In logical replication, this parameter also limits how often a failing replication apply worker will be respawned.

`recovery_min_apply_delay` (integer)

By default, a standby server restores WAL records from the sending server as soon as possible. It may be useful to have a time-delayed copy of the data, offering opportunities to correct data loss errors. This parameter allows you to delay recovery by a specified amount of time. For example, if you set this parameter to `5min`, the standby will replay each transaction commit only when the system time on the standby is at least five minutes past the commit time reported by the primary. If this value is specified without units, it is taken as milliseconds. The default is zero, adding no delay.

It is possible that the replication delay between servers exceeds the value of this parameter, in which case no delay is added. Note that the delay is calculated between the WAL time stamp as written on primary and the current time on the standby. Delays in transfer because of network lag or cascading replication configurations may reduce the actual wait time significantly. If the system clocks on primary and standby are not synchronized, this may lead to recovery applying records earlier than expected; but that is not a major issue because useful settings of this parameter are much larger than typical time deviations between servers.

The delay occurs only on WAL records for transaction commits. Other records are replayed as quickly as possible, which is not a problem because MVCC visibility rules ensure their effects are not visible until the corresponding commit record is applied.

The delay occurs once the database in recovery has reached a consistent state, until the standby is promoted or triggered. After that the standby will end recovery without further waiting.

WAL records must be kept on the standby until they are ready to be applied. Therefore, longer delays will result in a greater accumulation of WAL files, increasing disk space requirements for the standby's `pg_wal` directory.

This parameter is intended for use with streaming replication deployments; however, if the parameter is specified it will be honored in all cases except crash recovery. `hot_standby_feedback` will be delayed by use of this feature which could lead to bloat on the primary; use both together with care.

Warning

Synchronous replication is affected by this setting when `synchronous_commit` is set to `remote_apply`; every `COMMIT` will need to wait to be applied.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`sync_replication_slots` (boolean)

It enables a physical standby to synchronize logical failover slots from the primary server so that logical subscribers can resume replication from the new primary server after failover.

It is disabled by default. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.6.4. Subscribers

These settings control the behavior of a logical replication subscriber. Their values on the publisher are irrelevant. See [Section 29.11](#) for more details.

`max_replication_slots` (integer)

Specifies how many replication origins (see [Chapter 51](#)) can be tracked simultaneously, effectively limiting how many logical replication subscriptions can be created on the server. Setting it to a lower value than the current number of tracked replication origins (reflected in `pg_replication_origin_status`) will prevent the server from starting. `max_replication_slots` must be set to at least the number of subscriptions that will be added to the subscriber, plus some reserve for table synchronization.

Note that this parameter also applies on a sending server, but with a different meaning.

`max_logical_replication_workers` (integer)

Specifies maximum number of logical replication workers. This includes leader apply workers, parallel apply workers, and table synchronization workers.

For distributed system specifics, see [max_logical_replication_workers](#).

Logical replication workers are taken from the pool defined by `max_worker_processes`.

The default value is 4. This parameter can only be set at server start.

`max_sync_workers_per_subscription` (integer)

Maximum number of synchronization workers per subscription. This parameter controls the amount of parallelism of the initial data copy during the subscription initialization or when new tables are added.

Currently, there can be only one synchronization worker per table.

The synchronization workers are taken from the pool defined by `max_logical_replication_workers`.

The default value is 2. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`max_parallel_apply_workers_per_subscription` (integer)

Maximum number of parallel apply workers per subscription. This parameter controls the amount of parallelism for streaming of in-progress transactions with subscription parameter `streaming = parallel`.

The parallel apply workers are taken from the pool defined by `max_logical_replication_workers`.

The default value is 2. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.7. Query Planning

19.7.1. Planner Method Configuration

These configuration parameters provide a crude method of influencing the query plans chosen by the query optimizer. If the default plan chosen by the optimizer for a particular query is not optimal, a *temporary* solution is to use one of these configuration parameters to force the optimizer to choose a different plan. Better ways to improve the quality of the plans chosen by the optimizer include adjusting the planner cost constants (see [Section 19.7.2](#)), running `ANALYZE` manually, increasing the value of the [default_statistics_target](#) configuration parameter, and increasing the amount of statistics collected for specific columns using `ALTER TABLE SET STATISTICS`.

`enable_any_to_lateral_transformation` (boolean)

Enables or disables transformation of `ANY` subqueries into `LATERAL` joins. The default is `on`.

`enable_async_append` (boolean)

Enables or disables the query planner's use of async-aware append plan types. The default is `on`.

`enable_bitmapscan` (boolean)

Enables or disables the query planner's use of bitmap-scan plan types. The default is `on`.

`enable_extra_transformations` (boolean)

Enables or disables additional query-tree transformations. In some cases, the query tree is collapsed by “pulling up” the subqueries. The default is `on`.

`enable_gathermerge` (boolean)

Enables or disables the query planner's use of gather merge plan types. The default is `on`.

`enable_group_by_reordering` (boolean)

Controls if the query planner will produce a plan which will provide `GROUP BY` keys sorted in the order of keys of a child node of the plan, such as an index scan. When disabled, the query planner will produce a plan with `GROUP BY` keys only sorted to match the `ORDER BY` clause, if any. When enabled, the planner will try to produce a more efficient plan. The default value is `on`.

`enable_hashagg` (boolean)

Enables or disables the query planner's use of hashed aggregation plan types. The default is `on`.

`enable_alternative_sorting_cost_model` (boolean)

Enables or disables the query planner's use of alternative model for estimating the cost of tuple sorting. If this parameter is disabled, the query planner will use the standard model. The default is `on`.

`enable_hashjoin` (boolean)

Enables or disables the query planner's use of hash-join plan types. The default is `on`.

`enable_incremental_sort` (boolean)

Enables or disables the query planner's use of incremental sort steps. The default is `on`.

`enable_indexscan` (boolean)

Enables or disables the query planner's use of index-scan and index-only-scan plan types. The default is `on`. Also see [enable_indexonlyscan](#).

`enable_indexonlyscan` (boolean)

Enables or disables the query planner's use of index-only-scan plan types (see [Section 11.9](#)). The default is `on`. The [enable_indexscan](#) setting must also be enabled to have the query planner consider index-only-scans.

`enable_material` (boolean)

Enables or disables the query planner's use of materialization. It is impossible to suppress materialization entirely, but turning this variable off prevents the planner from inserting materialize nodes except in cases where it is required for correctness. The default is `on`.

`enable_memoize` (boolean)

Enables or disables the query planner's use of memoize plans for caching results from parameterized scans inside nested-loop joins. This plan type allows scans to the underlying plans to be skipped when the results for the current parameters are already in the cache. Less commonly looked up results may be evicted from the cache when more space is required for new entries. The default is `on`.

`enable_mergejoin` (boolean)

Enables or disables the query planner's use of merge-join plan types. The default is `on`.

`enable_nestloop` (boolean)

Enables or disables the query planner's use of nested-loop join plans. It is impossible to suppress nested-loop joins entirely, but turning this variable off discourages the planner from using one if there are other methods available. The default is `on`.

`enable_parallel_append` (boolean)

Enables or disables the query planner's use of parallel-aware append plan types. The default is `on`.

`enable_parallel_hash` (boolean)

Enables or disables the query planner's use of hash-join plan types with parallel hash. Has no effect if hash-join plans are not also enabled. The default is `on`.

`enable_parallel_temptables` (boolean)

Enables or disables parallel workers in queries involving temporary table scans. May degrade performance for workloads with frequent temporary table updates. The default value is `off`. Note that this feature is currently experimental.

Set the `enable_parallel_temptables` parameter to `on` only when the [enable_temp_memory_catalog](#) parameter is disabled. The corresponding features are mutually exclusive.

`enable_partition_pruning` (boolean)

Enables or disables the query planner's ability to eliminate a partitioned table's partitions from query plans. This also controls the planner's ability to generate query plans which allow the query executor to remove (ignore) partitions during query execution. The default is `on`. See [Section 5.12.4](#) for details.

`enable_partitionwise_join` (boolean)

Enables or disables the query planner's use of partitionwise join, which allows a join between partitioned tables to be performed by joining the matching partitions. Partitionwise join currently applies only when the join conditions include all the partition keys, which must be of the same data type and have one-to-one matching sets of child partitions. With this setting enabled, the number of nodes whose memory usage is restricted by `work_mem` appearing in the final plan can increase linearly according to the number of partitions being scanned. This can result in a large increase in overall memory consumption during the execution of the query. Query planning also becomes significantly more expensive in terms of memory and CPU. The default value is `off`.

For distributed system specifics, see [enable_partitionwise_aggregate](#).

`enable_partitionwise_aggregate` (boolean)

Enables or disables the query planner's use of partitionwise grouping or aggregation, which allows grouping or aggregation on partitioned tables to be performed separately for each partition. If the `GROUP BY` clause does not include the partition keys, only partial aggregation can be performed on a per-partition basis, and finalization must be performed later. With this setting enabled, the number of nodes whose memory usage is restricted by `work_mem` appearing in the final plan can increase linearly according to the number of partitions being scanned. This can result in a large increase in overall memory consumption during the execution of the query. Query planning also becomes significantly more expensive in terms of memory and CPU. The default value is `off`.

For distributed system specifics, see [enable_partitionwise_aggregate](#).

`enable_presorted_aggregate` (boolean)

Controls if the query planner will produce a plan which will provide rows which are presorted in the order required for the query's `ORDER BY` / `DISTINCT` aggregate functions. When disabled, the query planner will produce a plan which will always require the executor to perform a sort before performing aggregation of each aggregate function containing an `ORDER BY` or `DISTINCT` clause. When enabled, the planner will try to produce a more efficient plan which provides input to the aggregate functions which is presorted in the order they require for aggregation. The default value is `on`.

`enable_compound_index_stats` (boolean)

Enables or disables use of compound indexes statistics for selectivity estimation.

Default: `on`

`enable_self_join_removal` (boolean)

Enables or disables the query planner optimization, which analyses the query tree and replaces self joins with semantically equivalent single scans. Takes into consideration only plain tables. Removing self joins based on unique column can significantly speed up queries without affecting the results. The default is `on`.

`enable_seqscan` (boolean)

Enables or disables the query planner's use of sequential scan plan types. It is impossible to suppress sequential scans entirely, but turning this variable off discourages the planner from using one if there are other methods available. The default is `on`.

`enable_sort` (boolean)

Enables or disables the query planner's use of explicit sort steps. It is impossible to suppress explicit sorts entirely, but turning this variable off discourages the planner from using one if there are other methods available. The default is `on`.

`enable_tidscan` (boolean)

Enables or disables the query planner's use of TID scan plan types. The default is `on`.

`self_join_search_limit` (integer)

Specifies the maximum size of a list of links from a query to the same table where self joins will be looked for. Such a list is created to analyze a possibility of self join removal. To find a self join, inter-relationships of all the elements in this list with all the other ones must be analyzed. The limitation on the list size aims to reduce the quickly growing complexity of this process. The default is `32`.

`planner_upper_limit_estimation` (boolean)

Enables or disables the query planner to overestimate the expected number of rows in statements that contain a comparison to an unknown constant. This can be useful, for instance, when the planner chooses a plan with a nested loop join instead of a more optimal plan. The default value is `off`.

19.7.2. Planner Cost Constants

The *cost* variables described in this section are measured on an arbitrary scale. Only their relative values matter, hence scaling them all up or down by the same factor will result in no change in the planner's choices. By default, these cost variables are based on the cost of sequential page fetches; that is, `seq_page_cost` is conventionally set to `1.0` and the other cost variables are set with reference to that. But you can use a different scale if you prefer, such as actual execution times in milliseconds on a particular machine.

Note

Unfortunately, there is no well-defined method for determining ideal values for the cost variables. They are best treated as averages over the entire mix of queries that a particular installation will receive. This means that changing them on the basis of just a few experiments is very risky.

`seq_page_cost` (floating point)

Sets the planner's estimate of the cost of a disk page fetch that is part of a series of sequential fetches. The default is `1.0`. This value can be overridden for tables and indexes in a particular tablespace by setting the tablespace parameter of the same name (see [ALTER TABLESPACE](#)).

`random_page_cost` (floating point)

Sets the planner's estimate of the cost of a non-sequentially-fetched disk page. The default is `4.0`. This value can be overridden for tables and indexes in a particular tablespace by setting the tablespace parameter of the same name (see [ALTER TABLESPACE](#)).

Reducing this value relative to `seq_page_cost` will cause the system to prefer index scans; raising it will make index scans look relatively more expensive. You can raise or lower both values together to change the importance of disk I/O costs relative to CPU costs, which are described by the following parameters.

Random access to mechanical disk storage is normally much more expensive than four times sequential access. However, a lower default is used (4.0) because the majority of random accesses to disk, such as indexed reads, are assumed to be in cache. The default value can be thought of as modeling random access as 40 times slower than sequential, while expecting 90% of random reads to be cached.

If you believe a 90% cache rate is an incorrect assumption for your workload, you can increase `random_page_cost` to better reflect the true cost of random storage reads. Correspondingly, if your data is likely to be completely in cache, such as when the database is smaller than the total server memory, decreasing `random_page_cost` can be appropriate. Storage that has a low random read cost relative to sequential, e.g., solid-state drives, might also be better modeled with a lower value for `random_page_cost`, e.g., 1.1.

Tip

Although the system will let you set `random_page_cost` to less than `seq_page_cost`, it is not physically sensible to do so. However, setting them equal makes sense if the database is entirely cached in RAM, since in that case there is no penalty for touching pages out of sequence. Also, in a heavily-cached database you should lower both values relative to the CPU parameters, since the cost of fetching a page already in RAM is much smaller than it would normally be.

`write_page_cost` (floating point)

Sets the planner's estimate of the cost of flushing temporary table pages to disk. The default is 5.0. Note that this parameter only works when [enable_parallel_temptables](#) is enabled.

`cpu_tuple_cost` (floating point)

Sets the planner's estimate of the cost of processing each row during a query. The default is 0.01.

`cpu_index_tuple_cost` (floating point)

Sets the planner's estimate of the cost of processing each index entry during an index scan. The default is 0.005.

`cpu_operator_cost` (floating point)

Sets the planner's estimate of the cost of processing each operator or function executed during a query. The default is 0.0025.

`parallel_setup_cost` (floating point)

Sets the planner's estimate of the cost of launching parallel worker processes. The default is 1000.

`parallel_tuple_cost` (floating point)

Sets the planner's estimate of the cost of transferring one tuple from a parallel worker process to another process. The default is 0.1.

`min_parallel_table_scan_size` (integer)

Sets the minimum amount of table data that must be scanned in order for a parallel scan to be considered. For a parallel sequential scan, the amount of table data scanned is always equal to the size of the table, but when indexes are used the amount of table data scanned will normally be less. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default is 8 megabytes (8MB).

`min_parallel_index_scan_size` (integer)

Sets the minimum amount of index data that must be scanned in order for a parallel scan to be considered. Note that a parallel index scan typically won't touch the entire index; it is the number of pages which the planner believes will actually be touched by the scan which is relevant. This parameter is also used to decide whether a particular index can participate in a parallel vacuum. See [VACUUM](#). If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default is 512 kilobytes (512kB).

`effective_cache_size` (integer)

Sets the planner's assumption about the effective size of the disk cache that is available to a single query. This is factored into estimates of the cost of using an index; a higher value makes it more likely index scans will be used, a lower value makes it more likely sequential scans will be used. When setting this parameter you should consider both Postgres Pro's shared buffers and the portion of the kernel's disk cache that will be used for Postgres Pro data files, though some data might exist in both places. Also, take into account the expected number of concurrent queries on different tables, since they will have to share the available space. This parameter has no effect on the size of shared memory allocated by Postgres Pro, nor does it reserve kernel disk cache; it is used only for estimation purposes. The system also does not assume data remains in the disk cache between queries. If this value is specified without units, it is taken as blocks, that is `BLCKSZ` bytes, typically 8kB. The default is 4 gigabytes (4GB). (If `BLCKSZ` is not 8kB, the default value scales proportionally to it.)

`jit_above_cost` (floating point)

Sets the query cost above which JIT compilation is activated, if enabled (see [Chapter 30](#)). Performing JIT costs planning time but can accelerate query execution. Setting this to `-1` disables JIT compilation. The default is 100000.

`jit_inline_above_cost` (floating point)

Sets the query cost above which JIT compilation attempts to inline functions and operators. Inlining adds planning time, but can improve execution speed. It is not meaningful to set this to less than `jit_above_cost`. Setting this to `-1` disables inlining. The default is 500000.

`jit_optimize_above_cost` (floating point)

Sets the query cost above which JIT compilation applies expensive optimizations. Such optimization adds planning time, but can improve execution speed. It is not meaningful to set this to less than `jit_above_cost`, and it is unlikely to be beneficial to set it to more than `jit_inline_above_cost`. Setting this to `-1` disables expensive optimizations. The default is 500000.

`generic_plan_fuzz_factor` (double)

Sets the plan cost calculation coefficient of the planner, which increases the probability that the generic or custom plan will be selected more often. By default, the value is set to 1, which means that the generic plan is preferred over the custom plan. The higher the value, the more likely it is that the custom plan will be selected automatically. If [plan cache mode](#) is set to `force_generic_plan`, the planner will always opt for the generic plan, regardless of the value of this configuration parameter.

19.7.3. Genetic Query Optimizer

The genetic query optimizer (GEQO) is an algorithm that does query planning using heuristic searching. This reduces planning time for complex queries (those joining many relations), at the cost of producing plans that are sometimes inferior to those found by the normal exhaustive-search algorithm. For more information see [Chapter 61](#).

`geqo` (boolean)

Enables or disables genetic query optimization. This is on by default. It is usually best not to turn it off in production; the `geqo_threshold` variable provides more granular control of GEQO.

`geqo_threshold (integer)`

Use genetic query optimization to plan queries with at least this many `FROM` items involved. (Note that a `FULL OUTER JOIN` construct counts as only one `FROM` item.) The default is 12. For simpler queries it is usually best to use the regular, exhaustive-search planner, but for queries with many tables the exhaustive search takes too long, often longer than the penalty of executing a suboptimal plan. Thus, a threshold on the size of the query is a convenient way to manage use of GEQO.

`geqo_effort (integer)`

Controls the trade-off between planning time and query plan quality in GEQO. This variable must be an integer in the range from 1 to 10. The default value is five. Larger values increase the time spent doing query planning, but also increase the likelihood that an efficient query plan will be chosen.

`geqo_effort` doesn't actually do anything directly; it is only used to compute the default values for the other variables that influence GEQO behavior (described below). If you prefer, you can set the other parameters by hand instead.

`geqo_pool_size (integer)`

Controls the pool size used by GEQO, that is the number of individuals in the genetic population. It must be at least two, and useful values are typically 100 to 1000. If it is set to zero (the default setting) then a suitable value is chosen based on `geqo_effort` and the number of tables in the query.

`geqo_generations (integer)`

Controls the number of generations used by GEQO, that is the number of iterations of the algorithm. It must be at least one, and useful values are in the same range as the pool size. If it is set to zero (the default setting) then a suitable value is chosen based on `geqo_pool_size`.

`geqo_selection_bias (floating point)`

Controls the selection bias used by GEQO. The selection bias is the selective pressure within the population. Values can be from 1.50 to 2.00; the latter is the default.

`geqo_seed (floating point)`

Controls the initial value of the random number generator used by GEQO to select random paths through the join order search space. The value can range from zero (the default) to one. Varying the value changes the set of join paths explored, and may result in a better or worse best path being found.

19.7.4. Other Planner Options

`autoprepare_for_protocol (enum)`

Sets the protocol used for submitting queries that could be autoprepared. The allowed values are `simple`, `extended`, and `all` (both simple and extended query protocols can be used).

`autoprepare_for_protocol` takes effect only if [autoprepare_threshold](#) is set.

`autoprepare_limit (integer)`

Specifies the maximal number of statements that can be autoprepared on a backend. If this parameter is set to zero, there is no limit. Note that an infinite number of prepared queries can slow down query execution and cause backend memory overflow. The default value is 100.

`autoprepare_limit` takes effect only if [autoprepare_threshold](#) is set.

`autoprepare_memory_limit (integer)`

Limits the amount of memory that can be allocated for autoprepared statements on a backend. If this parameter is set to zero (the default), there is no memory limit. Setting this parameter to a non-zero value can cause a slowdown since calculating memory used by prepared statements adds some

overhead. To avoid performance penalty, you can limit the number of autoprepared statements using the [autoprepere_limit](#) parameter.

`autoprepere_memory_limit` takes effect only if [autoprepere_threshold](#) is set.

`autoprepere_threshold (integer)`

Specifies the minimal number of times a statement should be executed before it is autoprepared. If set to zero (the default), disables the autoprepere mode. See [Section 14.5](#) for more information.

`default_statistics_target (integer)`

Sets the default statistics target for table columns without a column-specific target set via `ALTER TABLE SET STATISTICS`. Larger values increase the time needed to do `ANALYZE`, but might improve the quality of the planner's estimates. The default is 100. For more information on the use of statistics by the Postgres Pro query planner, refer to [Section 14.2](#).

`constraint_exclusion (enum)`

Controls the query planner's use of table constraints to optimize queries. The allowed values of `constraint_exclusion` are `on` (examine constraints for all tables), `off` (never examine constraints), and `partition` (examine constraints only for inheritance child tables and `UNION ALL` subqueries). `partition` is the default setting. It is often used with traditional inheritance trees to improve performance.

When this parameter allows it for a particular table, the planner compares query conditions with the table's `CHECK` constraints, and omits scanning tables for which the conditions contradict the constraints. For example:

```
CREATE TABLE parent(key integer, ...);
CREATE TABLE child1000(check (key between 1000 and 1999)) INHERITS(parent);
CREATE TABLE child2000(check (key between 2000 and 2999)) INHERITS(parent);
...
SELECT * FROM parent WHERE key = 2400;
```

With constraint exclusion enabled, this `SELECT` will not scan `child1000` at all, improving performance.

Currently, constraint exclusion is enabled by default only for cases that are often used to implement table partitioning via inheritance trees. Turning it on for all tables imposes extra planning overhead that is quite noticeable on simple queries, and most often will yield no benefit for simple queries. If you have no tables that are partitioned using traditional inheritance, you might prefer to turn it off entirely. (Note that the equivalent feature for partitioned tables is controlled by a separate parameter, [enable_partition_pruning](#).)

Refer to [Section 5.12.5](#) for more information on using constraint exclusion to implement partitioning.

`cursor_tuple_fraction (floating point)`

Sets the planner's estimate of the fraction of a cursor's rows that will be retrieved. The default is 0.1. Smaller values of this setting bias the planner towards using “fast start” plans for cursors, which will retrieve the first few rows quickly while perhaps taking a long time to fetch all rows. Larger values put more emphasis on the total estimated time. At the maximum setting of 1.0, cursors are planned exactly like regular queries, considering only the total estimated time and not how soon the first rows might be delivered.

`from_collapse_limit (integer)`

The planner will merge sub-queries into upper queries if the resulting `FROM` list would have no more than this many items. Smaller values reduce planning time but might yield inferior query plans. The default is eight. For more information see [Section 14.3](#).

Setting this value to [geqo_threshold](#) or more may trigger use of the GEQO planner, resulting in non-optimal plans. See [Section 19.7.3](#).

`jit` (boolean)

Determines whether JIT compilation may be used by Postgres Pro, if available (see [Chapter 30](#)). The default is `on`.

`join_collapse_limit` (integer)

The planner will rewrite explicit `JOIN` constructs (except `FULL JOINS`) into lists of `FROM` items whenever a list of no more than this many items would result. Smaller values reduce planning time but might yield inferior query plans.

By default, this variable is set the same as `from_collapse_limit`, which is appropriate for most uses. Setting it to 1 prevents any reordering of explicit `JOINS`. Thus, the explicit join order specified in the query will be the actual order in which the relations are joined. Because the query planner does not always choose the optimal join order, advanced users can elect to temporarily set this variable to 1, and then specify the join order they desire explicitly. For more information see [Section 14.3](#).

Setting this value to [geqo_threshold](#) or more may trigger use of the GEQO planner, resulting in non-optimal plans. See [Section 19.7.3](#).

`plan_cache_mode` (enum)

Prepared statements (either explicitly prepared or implicitly generated, for example by PL/pgSQL) can be executed using custom or generic plans. Custom plans are made afresh for each execution using its specific set of parameter values, while generic plans do not rely on the parameter values and can be re-used across executions. Thus, use of a generic plan saves planning time, but if the ideal plan depends strongly on the parameter values then a generic plan may be inefficient. The choice between these options is normally made automatically, but it can be overridden with `plan_cache_mode`. The allowed values are `auto` (the default), `force_custom_plan` and `force_generic_plan`. This setting is considered when a cached plan is to be executed, not when it is prepared. For more information see [PREPARE](#).

`recursive_worktable_factor` (floating point)

Sets the planner's estimate of the average size of the working table of a [recursive query](#), as a multiple of the estimated size of the initial non-recursive term of the query. This helps the planner choose the most appropriate method for joining the working table to the query's other tables. The default value is 10.0. A smaller value such as 1.0 can be helpful when the recursion has low “fan-out” from one step to the next, as for example in shortest-path queries. Graph analytics queries may benefit from larger-than-default values.

`enable_appendorpath` (boolean)

Enables the Append plan for `OR` clauses. This parameter adds one more strategy for the optimizer: the Append plan for expressions containing `OR` clauses. It is useful for applications with auto-generated queries.

`seq_scan_startup_cost_first_row` (boolean)

Adds an average cost of getting the first tuple to the startup cost of sequential scan plans. With this option enabled, index scans get higher priority over sequential scans. If this option is off (the default), the startup cost includes only preliminary work before starting the scan.

19.8. Error Reporting and Logging

19.8.1. Where to Log

`log_destination` (string)

Postgres Pro supports several methods for logging server messages, including `stderr`, `csvlog`, `json-log`, and `syslog`. On Windows, `eventlog` is also supported. Set this parameter to a list of desired log

destinations separated by commas. The default is to log to stderr only. This parameter can only be set in the `postgresql.conf` file or on the server command line.

If `csvlog` is included in `log_destination`, log entries are output in “comma separated value” (CSV) format, which is convenient for loading logs into programs. See [Section 19.8.4](#) for details. `logging_collector` must be enabled to generate CSV-format log output.

If `jsonlog` is included in `log_destination`, log entries are output in JSON format, which is convenient for loading logs into programs. See [Section 19.8.5](#) for details. `logging_collector` must be enabled to generate JSON-format log output.

When either `stderr`, `csvlog` or `jsonlog` are included, the file `current_logfiles` is created to record the location of the log file(s) currently in use by the logging collector and the associated logging destination. This provides a convenient way to find the logs currently in use by the instance. Here is an example of this file's content:

```
stderr log/postgresql.log
csvlog log/postgresql.csv
jsonlog log/postgresql.json
```

`current_logfiles` is recreated when a new log file is created as an effect of rotation, and when `log_destination` is reloaded. It is removed when none of `stderr`, `csvlog` or `jsonlog` are included in `log_destination`, and when the logging collector is disabled.

Note

On most Unix systems, you will need to alter the configuration of your system's syslog daemon in order to make use of the `syslog` option for `log_destination`. Postgres Pro can log to syslog facilities `LOCAL0` through `LOCAL7` (see [syslog facility](#)), but the default syslog configuration on most platforms will discard all such messages. You will need to add something like:

```
local0.*      /var/log/postgresql
```

to the syslog daemon's configuration file to make it work.

On Windows, when you use the `eventlog` option for `log_destination`, you should register an event source and its library with the operating system so that the Windows Event Viewer can display event log messages cleanly. See [Section 18.12](#) for details.

`logging_collector` (boolean)

This parameter enables the *logging collector*, which is a background process that captures log messages sent to stderr and redirects them into log files. This approach is often more useful than logging to syslog, since some types of messages might not appear in syslog output. (One common example is dynamic-linker failure messages; another is error messages produced by scripts such as `archive_command`.) This parameter can only be set at server start.

Note

It is possible to log to stderr without using the logging collector; the log messages will just go to wherever the server's stderr is directed. However, that method is only suitable for low log volumes, since it provides no convenient way to rotate log files. Also, on some platforms not using the logging collector can result in lost or garbled log output, because multiple processes writing concurrently to the same log file can overwrite each other's output.

Note

The logging collector is designed to never lose messages. This means that in case of extremely high load, server processes could be blocked while trying to send additional log messages

when the collector has fallen behind. In contrast, syslog prefers to drop messages if it cannot write them, which means it may fail to log some messages in such cases but it will not block the rest of the system.

`log_directory (string)`

When `logging_collector` is enabled, this parameter determines the directory in which log files will be created. It can be specified as an absolute path, or relative to the cluster data directory. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is `log`.

`log_filename (string)`

When `logging_collector` is enabled, this parameter sets the file names of the created log files. The value is treated as a `strftime` pattern, so %-escapes can be used to specify time-varying file names. (Note that if there are any time-zone-dependent %-escapes, the computation is done in the zone specified by [log_timezone](#).) The supported %-escapes are similar to those listed in the Open Group's [strftime](#) specification. Note that the system's `strftime` is not used directly, so platform-specific (nonstandard) extensions do not work. The default is `postgresql-%Y-%m-%d_%H%M%S.log`.

If you specify a file name without escapes, you should plan to use a log rotation utility to avoid eventually filling the entire disk. In releases prior to 8.4, if no % escapes were present, PostgreSQL would append the epoch of the new log file's creation time, but this is no longer the case.

If CSV-format output is enabled in `log_destination`, `.csv` will be appended to the timestamped log file name to create the file name for CSV-format output. (If `log_filename` ends in `.log`, the suffix is replaced instead.)

If JSON-format output is enabled in `log_destination`, `.json` will be appended to the timestamped log file name to create the file name for JSON-format output. (If `log_filename` ends in `.log`, the suffix is replaced instead.)

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_file_mode (integer)`

On Unix systems this parameter sets the permissions for log files when `logging_collector` is enabled. (On Microsoft Windows this parameter is ignored.) The parameter value is expected to be a numeric mode specified in the format accepted by the `chmod` and `umask` system calls. (To use the customary octal format the number must start with a 0 (zero).)

The default permissions are `0600`, meaning only the server owner can read or write the log files. The other commonly useful setting is `0640`, allowing members of the owner's group to read the files. Note however that to make use of such a setting, you'll need to alter [log_directory](#) to store the files somewhere outside the cluster data directory. In any case, it's unwise to make the log files world-readable, since they might contain sensitive data.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_rotation_age (integer)`

When `logging_collector` is enabled, this parameter determines the maximum amount of time to use an individual log file, after which a new log file will be created. If this value is specified without units, it is taken as minutes. The default is 24 hours. Set to zero to disable time-based creation of new log files. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_rotation_size (integer)`

When `logging_collector` is enabled, this parameter determines the maximum size of an individual log file. After this amount of data has been emitted into a log file, a new log file will be created. If this

value is specified without units, it is taken as kilobytes. The default is 10 megabytes. Set to zero to disable size-based creation of new log files. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_truncate_on_rotation` (boolean)

When `logging_collector` is enabled, this parameter will cause Postgres Pro to truncate (overwrite), rather than append to, any existing log file of the same name. However, truncation will occur only when a new file is being opened due to time-based rotation, not during server startup or size-based rotation. When off, pre-existing files will be appended to in all cases. For example, using this setting in combination with a `log_filename` like `postgresql-%H.log` would result in generating twenty-four hourly log files and then cyclically overwriting them. This parameter can only be set in the `postgresql.conf` file or on the server command line.

Example: To keep 7 days of logs, one log file per day named `server_log.Mon`, `server_log.Tue`, etc., and automatically overwrite last week's log with this week's log, set `log_filename` to `server_log.%a`, `log_truncate_on_rotation` to on, and `log_rotation_age` to 1440.

Example: To keep 24 hours of logs, one log file per hour, but also rotate sooner if the log file size exceeds 1GB, set `log_filename` to `server_log.%H%M`, `log_truncate_on_rotation` to on, `log_rotation_age` to 60, and `log_rotation_size` to 1000000. Including `%M` in `log_filename` allows any size-driven rotations that might occur to select a file name different from the hour's initial file name.

`syslog_facility` (enum)

When logging to syslog is enabled, this parameter determines the syslog “facility” to be used. You can choose from `LOCAL0`, `LOCAL1`, `LOCAL2`, `LOCAL3`, `LOCAL4`, `LOCAL5`, `LOCAL6`, `LOCAL7`; the default is `LOCAL0`. See also the documentation of your system's syslog daemon. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`syslog_ident` (string)

When logging to syslog is enabled, this parameter determines the program name used to identify Postgres Pro messages in syslog logs. The default is `postgres`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`syslog_sequence_numbers` (boolean)

When logging to syslog and this is on (the default), then each message will be prefixed by an increasing sequence number (such as [2]). This circumvents the “--- last message repeated N times ---” suppression that many syslog implementations perform by default. In more modern syslog implementations, repeated message suppression can be configured (for example, `$RepeatedMsgReduction` in `rsyslog`), so this might not be necessary. Also, you could turn this off if you actually want to suppress repeated messages.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`syslog_split_messages` (boolean)

When logging to syslog is enabled, this parameter determines how messages are delivered to syslog. When on (the default), messages are split by lines, and long lines are split so that they will fit into 1024 bytes, which is a typical size limit for traditional syslog implementations. When off, Postgres Pro server log messages are delivered to the syslog service as is, and it is up to the syslog service to cope with the potentially bulky messages.

If syslog is ultimately logging to a text file, then the effect will be the same either way, and it is best to leave the setting on, since most syslog implementations either cannot handle large messages or would need to be specially configured to handle them. But if syslog is ultimately writing into some other medium, it might be necessary or more useful to keep messages logically together.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`event_source (string)`

When logging to event log is enabled, this parameter determines the program name used to identify Postgres Pro messages in the log. The default is `Postgres Pro`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.8.2. When to Log

`log_min_messages (enum)`

Controls which [message levels](#) are written to the server log. Valid values are `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `INFO`, `NOTICE`, `WARNING`, `ERROR`, `LOG`, `FATAL`, and `PANIC`. Each level includes all the levels that follow it. The later the level, the fewer messages are sent to the log. The default is `WARNING`. Note that `LOG` has a different rank here than in [client_min_messages](#). Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_min_error_statement (enum)`

Controls which SQL statements that cause an error condition are recorded in the server log. The current SQL statement is included in the log entry for any message of the specified [severity](#) or higher. Valid values are `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `INFO`, `NOTICE`, `WARNING`, `ERROR`, `LOG`, `FATAL`, and `PANIC`. The default is `ERROR`, which means statements causing errors, log messages, fatal errors, or panics will be logged. To effectively turn off logging of failing statements, set this parameter to `PANIC`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_min_duration_statement (integer)`

Causes the duration of each completed statement to be logged if the statement ran for at least the specified amount of time. For example, if you set it to `250ms` then all SQL statements that run 250ms or longer will be logged. Enabling this parameter can be helpful in tracking down unoptimized queries in your applications. If this value is specified without units, it is taken as milliseconds. Setting this to zero prints all statement durations. `-1` (the default) disables logging statement durations. Only superusers and users with the appropriate `SET` privilege can change this setting.

This overrides [log_min_duration_sample](#), meaning that queries with duration exceeding this setting are not subject to sampling and are always logged.

For clients using extended query protocol, durations of the Parse, Bind, and Execute steps are logged independently.

Note

When using this option together with [log_statement](#), the text of statements that are logged because of `log_statement` will not be repeated in the duration log message. If you are not using syslog, it is recommended that you log the PID or session ID using [log_line_prefix](#) so that you can link the statement message to the later duration message using the process ID or session ID.

`log_min_duration_sample (integer)`

Allows sampling the duration of completed statements that ran for at least the specified amount of time. This produces the same kind of log entries as [log_min_duration_statement](#), but only for a subset of the executed statements, with sample rate controlled by [log_statement_sample_rate](#). For example, if you set it to `100ms` then all SQL statements that run 100ms or longer will be considered for sampling. Enabling this parameter can be helpful when the traffic is too high to log all queries. If this value is specified without units, it is taken as milliseconds. Setting this to zero samples all statement durations. `-1` (the default) disables sampling statement durations. Only superusers and users with the appropriate `SET` privilege can change this setting.

This setting has lower priority than `log_min_duration_statement`, meaning that statements with durations exceeding `log_min_duration_statement` are not subject to sampling and are always logged.

Other notes for `log_min_duration_statement` apply also to this setting.

`log_statement_sample_rate` (floating point)

Determines the fraction of statements with duration exceeding `log_min_duration_sample` that will be logged. Sampling is stochastic, for example `0.5` means there is statistically one chance in two that any given statement will be logged. The default is `1.0`, meaning to log all sampled statements. Setting this to zero disables sampled statement-duration logging, the same as setting `log_min_duration_sample` to `-1`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_transaction_sample_rate` (floating point)

Sets the fraction of transactions whose statements are all logged, in addition to statements logged for other reasons. It applies to each new transaction regardless of its statements' durations. Sampling is stochastic, for example `0.1` means there is statistically one chance in ten that any given transaction will be logged. `log_transaction_sample_rate` can be helpful to construct a sample of transactions. The default is `0`, meaning not to log statements from any additional transactions. Setting this to `1` logs all statements of all transactions. Only superusers and users with the appropriate `SET` privilege can change this setting.

Note

Like all statement-logging options, this option can add significant overhead.

`log_startup_progress_interval` (integer)

Sets the amount of time after which the startup process will log a message about a long-running operation that is still in progress, as well as the interval between further progress messages for that operation. The default is 10 seconds. A setting of `0` disables the feature. If this value is specified without units, it is taken as milliseconds. This setting is applied separately to each operation. This parameter can only be set in the `postgresql.conf` file or on the server command line.

For example, if syncing the data directory takes 25 seconds and thereafter resetting unlogged relations takes 8 seconds, and if this setting has the default value of 10 seconds, then a messages will be logged for syncing the data directory after it has been in progress for 10 seconds and again after it has been in progress for 20 seconds, but nothing will be logged for resetting unlogged relations.

[Table 19.2](#) explains the message severity levels used by Postgres Pro. If logging output is sent to syslog or Windows' eventlog, the severity levels are translated as shown in the table.

Table 19.2. Message Severity Levels

Severity	Usage	syslog	eventlog
DEBUG1 .. DEBUG5	Provides successively-more-detailed information for use by developers.	DEBUG	INFORMATION
INFO	Provides information implicitly requested by the user, e.g., output from <code>VACUUM VERBOSE</code> .	INFO	INFORMATION
NOTICE	Provides information that might be helpful to users, e.g., notice of truncation of long identifiers.	NOTICE	INFORMATION
WARNING	Provides warnings of likely problems, e.g., <code>COMMIT</code> outside a transaction block.	NOTICE	WARNING

Severity	Usage	syslog	eventlog
ERROR	Reports an error that caused the current command to abort.	WARNING	ERROR
LOG	Reports information of interest to administrators, e.g., checkpoint activity.	INFO	INFORMATION
FATAL	Reports an error that caused the current session to abort.	ERR	ERROR
PANIC	Reports an error that caused all database sessions to abort.	CRIT	ERROR

19.8.3. What to Log

Note

What you choose to log can have security implications; see [Section 24.4](#).

`application_name` (string)

The `application_name` can be any string of less than `NAMEDATALEN` characters (64 characters in a standard build). It is typically set by an application upon connection to the server. The name will be displayed in the `pg_stat_activity` view and included in CSV log entries. It can also be included in regular log entries via the `log_line_prefix` parameter. Only printable ASCII characters may be used in the `application_name` value. Other characters are replaced with [C-style hexadecimal escapes](#).

`debug_print_parse` (boolean)

`debug_print_rewritten` (boolean)

`debug_print_plan` (boolean)

These parameters enable various debugging output to be emitted. When set, they print the resulting parse tree, the query rewriter output, or the execution plan for each executed query. These messages are emitted at `LOG` message level, so by default they will appear in the server log but will not be sent to the client. You can change that by adjusting [client_min_messages](#) and/or [log_min_messages](#). These parameters are off by default.

`debug_pretty_print` (boolean)

When set, `debug_pretty_print` indents the messages produced by `debug_print_parse`, `debug_print_rewritten`, or `debug_print_plan`. This results in more readable but much longer output than the “compact” format used when it is off. It is on by default.

`log_autovacuum_min_duration` (integer)

Causes each action executed by autovacuum to be logged if it ran for at least the specified amount of time. Setting this to zero logs all autovacuum actions. `-1` disables logging autovacuum actions. If this value is specified without units, it is taken as milliseconds. For example, if you set this to `250ms` then all automatic vacuums and analyzes that run 250ms or longer will be logged. In addition, when this parameter is set to any value other than `-1`, a message will be logged if an autovacuum action is skipped due to a conflicting lock or a concurrently dropped relation. The default is `10min`. Enabling this parameter can be helpful in tracking autovacuum activity. This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`log_checkpoints` (boolean)

Causes checkpoints and restartpoints to be logged in the server log. Some statistics are included in the log messages, including the number of buffers written and the time spent writing them. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is on.

`log_connections` (boolean)

Causes each attempted connection to the server to be logged, as well as successful completion of both client authentication (if necessary) and authorization. Only superusers and users with the appropriate `SET` privilege can change this parameter at session start, and it cannot be changed at all within a session. The default is `off`.

Note

Some client programs, like `psql`, attempt to connect twice while determining if a password is required, so duplicate “connection received” messages do not necessarily indicate a problem.

`log_disconnections` (boolean)

Causes session terminations to be logged. The log output provides information similar to `log_connections`, plus the duration of the session. Only superusers and users with the appropriate `SET` privilege can change this parameter at session start, and it cannot be changed at all within a session. The default is `off`.

`log_duration` (boolean)

Causes the duration of every completed statement to be logged. The default is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

For clients using extended query protocol, durations of the Parse, Bind, and Execute steps are logged independently.

Note

The difference between enabling `log_duration` and setting `log_min_duration_statement` to zero is that exceeding `log_min_duration_statement` forces the text of the query to be logged, but this option doesn't. Thus, if `log_duration` is on and `log_min_duration_statement` has a positive value, all durations are logged but the query text is included only for statements exceeding the threshold. This behavior can be useful for gathering statistics in high-load installations.

`log_error_verbosity` (enum)

Controls the amount of detail written in the server log for each message that is logged. Valid values are `TERSE`, `DEFAULT`, and `VERBOSE`, each adding more fields to displayed messages. `TERSE` excludes the logging of `DETAIL`, `HINT`, `QUERY`, and `CONTEXT` error information. `VERBOSE` output includes the `SQLSTATE` error code (see also [Appendix A](#)) and the source code file name, function name, and line number that generated the error. Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_hostname` (boolean)

By default, connection log messages only show the IP address of the connecting host. Turning this parameter on causes logging of the host name as well. Note that depending on your host name resolution setup this might impose a non-negligible performance penalty. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_line_prefix` (string)

This is a `printf`-style string that is output at the beginning of each log line. `%` characters begin “escape sequences” that are replaced with status information as outlined below. Unrecognized escapes are ignored. Other characters are copied straight to the log line. Some escapes are only recognized by session processes, and will be treated as empty by background processes such as the main server process. Status information may be aligned either left or right by specifying a numeric literal after

the % and before the option. A negative value will cause the status information to be padded on the right with spaces to give it a minimum width, whereas a positive value will pad on the left. Padding can be useful to aid human readability in log files.

This parameter can only be set in the `postgresql.conf` file or on the server command line. The default is '`%m [%p]`' which logs a time stamp and the process ID.

Escape	Effect	Session only
%a	Application name	yes
%u	User name	yes
%d	Database name	yes
%r	Remote host name or IP address, and remote port	yes
%h	Remote host name or IP address	yes
%b	Backend type	no
%p	Process ID	no
%P	Process ID of the parallel group leader, if this process is a parallel query worker	no
%t	Time stamp without milliseconds	no
%m	Time stamp with milliseconds	no
%n	Time stamp with milliseconds (as a Unix epoch)	no
%i	Command tag: type of session's current command	yes
%e	SQLSTATE error code	no
%c	Session ID: see below	no
%l	Number of the log line for each session or process, starting at 1	no
%s	Process start time stamp	no
%v	Virtual transaction ID (proc-Number/localXID); see Section 67.1	no
%x	Transaction ID (0 if none is assigned); see Section 67.1	no
%q	Produces no output, but tells non-session processes to stop at this point in the string; ignored by session processes	no
%Q	Query identifier of the current query. Query identifiers are not computed by default, so this field will be zero unless compute_query_id parameter is enabled or a third-party module that computes query identifiers is configured.	yes
%%	Literal %	no

The `backend` type corresponds to the column `backend_type` in the view `pg_stat_activity`, but additional types can appear in the log that don't show in that view.

The `%c` escape prints a quasi-unique session identifier, consisting of two 4-byte hexadecimal numbers (without leading zeros) separated by a dot. The numbers are the process start time and the process ID, so `%c` can also be used as a space saving way of printing those items. For example, to generate the session identifier from `pg_stat_activity`, use this query:

```
SELECT to_hex(trunc(EXTRACT(EPOCH FROM backend_start)::integer) || '.' ||
        to_hex(pid)
FROM pg_stat_activity;
```

Tip

If you set a nonempty value for `log_line_prefix`, you should usually make its last character be a space, to provide visual separation from the rest of the log line. A punctuation character can be used too.

Tip

Syslog produces its own time stamp and process ID information, so you probably do not want to include those escapes if you are logging to syslog.

Tip

The `%q` escape is useful when including information that is only available in session (backend) context like user or database name. For example:

```
log_line_prefix = '%m [%p] %q%u@%d/%a '
```

Note

The `%Q` escape always reports a zero identifier for lines output by `log_statement` because `log_statement` generates output before an identifier can be calculated, including invalid statements for which an identifier cannot be calculated.

`log_lock_waits` (boolean)

Controls whether a log message is produced when a session waits longer than `deadlock_timeout` to acquire a lock. This is useful in determining if lock waits are causing poor performance. The default is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_recovery_conflict_waits` (boolean)

Controls whether a log message is produced when the startup process waits longer than `deadlock_timeout` for recovery conflicts. This is useful in determining if recovery conflicts prevent the recovery from applying WAL.

The default is `off`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log_parameter_max_length` (integer)

If greater than zero, each bind parameter value logged with a non-error statement-logging message is trimmed to this many bytes. Zero disables logging of bind parameters for non-error statement logs.

-1 (the default) allows bind parameters to be logged in full. If this value is specified without units, it is taken as bytes. Only superusers and users with the appropriate `SET` privilege can change this setting.

This setting only affects log messages printed as a result of [log_statement](#), [log_duration](#), and related settings. Non-zero values of this setting add some overhead, particularly if parameters are sent in binary form, since then conversion to text is required.

`log_parameter_max_length_on_error` (integer)

If greater than zero, each bind parameter value reported in error messages is trimmed to this many bytes. Zero (the default) disables including bind parameters in error messages. -1 allows bind parameters to be printed in full. If this value is specified without units, it is taken as bytes.

Non-zero values of this setting add overhead, as Postgres Pro will need to store textual representations of parameter values in memory at the start of each statement, whether or not an error eventually occurs. The overhead is greater when bind parameters are sent in binary form than when they are sent as text, since the former case requires data conversion while the latter only requires copying the string.

`log_statement` (enum)

Controls which SQL statements are logged. Valid values are `none` (off), `ddl`, `mod`, and `all` (all statements). `ddl` logs all data definition statements, such as `CREATE`, `ALTER`, and `DROP` statements. `mod` logs all `ddl` statements, plus data-modifying statements such as `INSERT`, `UPDATE`, `DELETE`, `TRUNCATE`, and `COPY FROM`. `PREPARE`, `EXECUTE`, and `EXPLAIN ANALYZE` statements are also logged if their contained command is of an appropriate type. For clients using extended query protocol, logging occurs when an `Execute` message is received, and values of the Bind parameters are included (with any embedded single-quote marks doubled).

The default is `none`. Only superusers and users with the appropriate `SET` privilege can change this setting.

Note

Statements that contain simple syntax errors are not logged even by the `log_statement = all` setting, because the log message is emitted only after basic parsing has been done to determine the statement type. In the case of extended query protocol, this setting likewise does not log statements that fail before the `Execute` phase (i.e., during parse analysis or planning). Set `log_min_error_statement` to `ERROR` (or lower) to log such statements.

Logged statements might reveal sensitive data and even contain plaintext passwords.

`log_next_xid_assign_threshold` (integer 64)

If logging is enabled by this setting, a log entry is emitted when during WAL replay the next XID is greater than the current XID by the specified value. A value of zero disables such logging, the default value is 100000000. Only superusers can change this setting.

`log_replication_commands` (boolean)

Causes each replication command and `walsender` process's replication slot acquisition/release to be logged in the server log. See [Section 56.4](#) for more information about replication command. The default value is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`shardman.context_log` (bool)

Logs the remote contexts. If enabled, in case of an error, displays a field `Remote CONTEXT`. Note that if the standart log level is set to `log_verbosity=terse`, the `shardman.context_log` will be disabled automatically.

Default: on.

`log_temp_files` (integer)

Controls logging of temporary file names and sizes. Temporary files can be created for sorts, hashes, and temporary query results. If enabled by this setting, a log entry is emitted for each temporary file, with the file size specified in bytes, when it is deleted. A value of zero logs all temporary file information, while positive values log only files whose size is greater than or equal to the specified amount of data. If this value is specified without units, it is taken as kilobytes. The default setting is -1, which disables such logging. Only superusers and users with the appropriate `SET` privilege can change this setting.

`log_timezone` (string)

Sets the time zone used for timestamps written in the server log. Unlike [TimeZone](#), this value is cluster-wide, so that all sessions will report timestamps consistently. The built-in default is `GMT`, but that is typically overridden in `postgresql.conf`; `initdb` will install a setting there corresponding to its system environment. See [Section 8.5.3](#) for more information. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.8.4. Using CSV-Format Log Output

Including `csvlog` in the `log_destination` list provides a convenient way to import log files into a database table. This option emits log lines in comma-separated-values (CSV) format, with these columns: time stamp with milliseconds, user name, database name, process ID, client host:port number, session ID, per-session line number, command tag, session start time, virtual transaction ID, regular transaction ID, error severity, SQLSTATE code, error message, error message detail, hint, internal query that led to the error (if any), character count of the error position therein, error context, user query that led to the error (if any and enabled by `log_min_error_statement`), character count of the error position therein, location of the error in the Postgres Pro source code (if `log_error_verbosity` is set to `verbose`), application name, backend type, process ID of parallel group leader, and query id. Here is a sample table definition for storing CSV-format log output:

```
CREATE TABLE postgres_log
(
    log_time timestamp(3) with time zone,
    user_name text,
    database_name text,
    process_id integer,
    connection_from text,
    session_id text,
    session_line_num bigint,
    command_tag text,
    session_start_time timestamp with time zone,
    virtual_transaction_id text,
    transaction_id bigint,
    error_severity text,
    sql_state_code text,
    message text,
    detail text,
    hint text,
    internal_query text,
    internal_query_pos integer,
    context text,
    query text,
    query_pos integer,
    location text,
    application_name text,
    backend_type text,
    leader_pid integer,
```

```

    query_id bigint,
    PRIMARY KEY (session_id, session_line_num)
);

```

To import a log file into this table, use the `COPY FROM` command:

```
COPY postgres_log FROM '/full/path/to/logfile.csv' WITH csv;
```

There are a few things you need to do to simplify importing CSV log files:

1. Set `log_filename` and `log_rotation_age` to provide a consistent, predictable naming scheme for your log files. This lets you predict what the file name will be and know when an individual log file is complete and therefore ready to be imported.
2. Set `log_rotation_size` to 0 to disable size-based log rotation, as it makes the log file name difficult to predict.
3. Set `log_truncate_on_rotation` to `on` so that old log data isn't mixed with the new in the same file.
4. The table definition above includes a primary key specification. This is useful to protect against accidentally importing the same information twice. The `COPY` command commits all of the data it imports at one time, so any error will cause the entire import to fail. If you import a partial log file and later import the file again when it is complete, the primary key violation will cause the import to fail. Wait until the log is complete and closed before importing. This procedure will also protect against accidentally importing a partial line that hasn't been completely written, which would also cause `COPY` to fail.

19.8.5. Using JSON-Format Log Output

Including `jsonlog` in the `log_destination` list provides a convenient way to import log files into many different programs. This option emits log lines in JSON format.

String fields with null values are excluded from output. Additional fields may be added in the future. User applications that process `jsonlog` output should ignore unknown fields.

Each log line is serialized as a JSON object with the set of keys and their associated values shown in [Table 19.3](#).

Table 19.3. Keys and Values of JSON Log Entries

Key name	Type	Description
timestamp	string	Time stamp with milliseconds
user	string	User name
dbname	string	Database name
pid	number	Process ID
remote_host	string	Client host
remote_port	number	Client port
session_id	string	Session ID
line_num	number	Per-session line number
ps	string	Current ps display
session_start	string	Session start time
vxid	string	Virtual transaction ID
txid	string	Regular transaction ID
error_severity	string	Error severity
state_code	string	SQLSTATE code
message	string	Error message

Key name	Type	Description
detail	string	Error message detail
hint	string	Error message hint
internal_query	string	Internal query that led to the error
internal_position	number	Cursor index into internal query
context	string	Error context
statement	string	Client-supplied query string
cursor_position	number	Cursor index into query string
func_name	string	Error location function name
file_name	string	File name of error location
file_line_num	number	File line number of the error location
application_name	string	Client application name
backend_type	string	Type of backend
leader_pid	number	Process ID of leader for active parallel workers
query_id	number	Query ID

19.8.6. Process Title

These settings control how process titles of server processes are modified. Process titles are typically viewed using programs like `ps` or, on Windows, Process Explorer. See [Section 27.1](#) for details.

`cluster_name` (string)

Sets a name that identifies this database cluster (instance) for various purposes. The cluster name appears in the process title for all server processes in this cluster. Moreover, it is the default application name for a standby connection (see [synchronous_standby_names](#)).

The name can be any string of less than `NAMEDATALEN` characters (64 characters in a standard build). Only printable ASCII characters may be used in the `cluster_name` value. Other characters are replaced with [C-style hexadecimal escapes](#). No name is shown if this parameter is set to the empty string `''` (which is the default). This parameter can only be set at server start.

`update_process_title` (boolean)

Enables updating of the process title every time a new SQL command is received by the server. This setting defaults to `on` on most platforms, but it defaults to `off` on Windows due to that platform's larger overhead for updating the process title. Only superusers and users with the appropriate `SET` privilege can change this setting.

19.8.7. PostgreSQL Logs

Postgres Pro Shardman uses standard PostgreSQL logging settings. They should be specified in [sdm-spec.json](#) in the `pgParameters` section, as shown in the example below:

```
{
  "ShardSpec": {
    "pgParameters": {
      "logging_collector": "on",
      ...
    },
    ...
  }
}
```



```
    },  
    ...  
}
```

By default, the following values of logging settings are specified in the configuration for all new Postgres Pro Shardman installations:

```
"logging_collector": "on",  
"log_directory": "{{dataDir}}",  
"log_filename": "postgres-{{cluster}}-{{shard}}-{{host}}-{{keeperID}}.log",  
"log_line_prefix": = "%m [%r] [%p]"
```

According to these settings, logs are written to the following file: `/var/lib/pgpro/sdm-VERSION/data/postgres-CLUSTERNAME-shard-N-HOSTNAME-keeper_M.log`. Here, `VERSION` is the major Postgres Pro Shardman version, `CLUSTERNAME` is the name of the current cluster, `shard-N` is the name of the current shard, `HOSTNAME` is the current host name, `keeper_M` is the ID of the integrated keeper process.

19.8.8. shardmand Logs

shardmand is a systemd unit, its logs are written to journald. You can use `journalctl` to examine it. For example, you can use the following command:

```
$ journalctl -u shardmand@cluster0.service
```

You can filter logs by arbitrary time limits using the `--since` and `--until` options, which restrict the entries displayed to those after or before the given time, respectively. The time values can come in a variety of formats. For absolute time values, you should use `YYYY-MM-DD HH:MM:SS`. For instance, we can see all of the entries since January 10th, 2023 at 5:15 PM by typing:

```
$ journalctl -u shardmand@cluster0.service --since "2023-01-10  
17:15:00"
```

If components of the above format are left off, some defaults will be applied. For instance, if the date is omitted, the current date will be assumed. If the time component is missing, "00:00:00" (midnight) will be substituted. The seconds field can be left off as well to default to "00":

```
$ journalctl -u shardmand@cluster0.service --since "2023-01-10" --  
until "2023-01-11 03:00"
```

To control the log verbosity for all Postgres Pro Shardman services, set `SDM_LOG_LEVEL` in the `shardmand` configuration file.

19.9. Run-time Statistics

19.9.1. Cumulative Query and Index Statistics

These parameters control the server-wide cumulative statistics system. When enabled, the data that is collected can be accessed via the `pg_stat` and `pg_statio` family of system views. Refer to [Chapter 27](#) for more information.

`track_activities` (boolean)

Enables the collection of information on the currently executing command of each session, along with its identifier and the time when that command began execution. This parameter is on by default. Note that even when enabled, this information is only visible to superusers, roles with privileges of the `pg_read_all_stats` role and the user owning the sessions being reported on (including sessions belonging to a role they have the privileges of), so it should not represent a security risk. Only superusers and users with the appropriate `SET` privilege can change this setting.

`track_activity_query_size` (integer)

Specifies the amount of memory reserved to store the text of the currently executing command for each active session, for the `pg_stat_activity.query` field. If this value is specified without units, it is taken as bytes. The default value is 1024 bytes. This parameter can only be set at server start.

`track_counts` (boolean)

Enables collection of statistics on database activity. This parameter is on by default, because the autovacuum daemon needs the collected information. Only superusers and users with the appropriate `SET` privilege can change this setting.

`track_io_timing` (boolean)

Enables timing of database I/O calls. This parameter is off by default, as it will repeatedly query the operating system for the current time, which may cause significant overhead on some platforms. You can use the `pg_test_timing` tool to measure the overhead of timing on your system. I/O timing information is displayed in `pg_stat_database`, `pg_stat_io`, in the output of `EXPLAIN` when the `BUFFERS` option is used, in the output of `VACUUM` when the `VERBOSE` option is used, by autovacuum for auto-vacuums and auto-analyzes, when `log_autovacuum_min_duration` is set and by `pg_stat_statements`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`track_wal_io_timing` (boolean)

Enables timing of WAL I/O calls. This parameter is off by default, as it will repeatedly query the operating system for the current time, which may cause significant overhead on some platforms. You can use the `pg_test_timing` tool to measure the overhead of timing on your system. I/O timing information is displayed in `pg_stat_wal`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`track_functions` (enum)

Enables tracking of function call counts and time used. Specify `pl` to track only procedural-language functions, `all` to also track SQL and C language functions. The default is `none`, which disables function statistics tracking. Only superusers and users with the appropriate `SET` privilege can change this setting.

Note

SQL-language functions that are simple enough to be “inlined” into the calling query will not be tracked, regardless of this setting.

`stats_fetch_consistency` (enum)

Determines the behavior when cumulative statistics are accessed multiple times within a transaction. When set to `none`, each access re-fetches counters from shared memory. When set to `cache`, the first access to statistics for an object caches those statistics until the end of the transaction unless `pg_stat_clear_snapshot()` is called. When set to `snapshot`, the first statistics access caches all statistics accessible in the current database, until the end of the transaction unless `pg_stat_clear_snapshot()` is called. Changing this parameter in a transaction discards the statistics snapshot. The default is `cache`.

Note

`none` is most suitable for monitoring systems. If values are only accessed once, it is the most efficient. `cache` ensures repeat accesses yield the same values, which is important for queries involving e.g. self-joins. `snapshot` can be useful when interactively inspecting statistics, but has higher overhead, particularly if many database objects exist.

19.9.2. Statistics Monitoring

`compute_query_id` (enum)

Enables in-core computation of a query identifier. Query identifiers can be displayed in the `pg_stat_activity` view, using `EXPLAIN`, or emitted in the log if configured via the `log_line_prefix` parameter. The `pg_stat_statements` extension also requires a query identifier to be computed. Note that an external module can alternatively be used if the in-core query identifier computation method is not acceptable. In this case, in-core computation must be always disabled. Valid values are `off` (always disabled), `on` (always enabled), `auto`, which lets modules such as `pg_stat_statements` automatically enable it, and `regress` which has the same effect as `auto`, except that the query identifier is not shown in the `EXPLAIN` output in order to facilitate automated regression testing. The default is `auto`.

Note

To ensure that only one query identifier is calculated and displayed, extensions that calculate query identifiers should throw an error if a query identifier has already been computed.

`log_statement_stats` (boolean)

`log_parser_stats` (boolean)

`log_planner_stats` (boolean)

`log_executor_stats` (boolean)

For each query, output performance statistics of the respective module to the server log. This is a crude profiling instrument, similar to the Unix `getrusage()` operating system facility. `log_statement_stats` reports total statement statistics, while the others report per-module statistics. `log_statement_stats` cannot be enabled together with any of the per-module options. All of these options are disabled by default. Only superusers and users with the appropriate `SET` privilege can change these settings.

19.10. Automatic Vacuuming

These settings control the behavior of the *autovacuum* feature. Refer to [Section 24.1.6](#) for more information. Note that many of these settings can be overridden on a per-table basis; see [Storage Parameters](#).

`autovacuum` (boolean)

Controls whether the server should run the autovacuum launcher daemon. This is on by default; however, `track_counts` must also be enabled for autovacuum to work. This parameter can only be set in the `postgresql.conf` file or on the server command line; however, autovacuuming can be disabled for individual tables by changing table storage parameters.

Note that even when this parameter is disabled, the system will launch autovacuum processes if necessary to shrink `pg_xact` and `pg_multixact`. See [Section 24.1.5](#) for more information.

`autovacuum_max_workers` (integer)

Specifies the maximum number of autovacuum processes (other than the autovacuum launcher) that may be running at any one time. The default is three. This parameter can only be set at server start.

`autovacuum_naptime` (integer)

Specifies the minimum delay between autovacuum runs on any given database. In each round the daemon examines the database and issues `VACUUM` and `ANALYZE` commands as needed for tables in that database. If this value is specified without units, it is taken as seconds. The default is one minute (1min). This parameter can only be set in the `postgresql.conf` file or on the server command line.

`autovacuum_vacuum_threshold` (integer)

Specifies the minimum number of updated or deleted tuples needed to trigger a `VACUUM` in any one table. The default is 50 tuples. This parameter can only be set in the `postgresql.conf` file or on

the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_vacuum_insert_threshold (integer)`

Specifies the number of inserted tuples needed to trigger a `VACUUM` in any one table. The default is 1000 tuples. If -1 is specified, autovacuum will not trigger a `VACUUM` operation on any tables based on the number of inserts. This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_analyze_threshold (integer)`

Specifies the minimum number of inserted, updated or deleted tuples needed to trigger an `ANALYZE` in any one table. The default is 50 tuples. This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_vacuum_scale_factor (floating point)`

Specifies a fraction of the table size to add to `autovacuum_vacuum_threshold` when deciding whether to trigger a `VACUUM`. The default is 0.2 (20% of table size). This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_vacuum_insert_scale_factor (floating point)`

Specifies a fraction of the table size to add to `autovacuum_vacuum_insert_threshold` when deciding whether to trigger a `VACUUM`. The default is 0.2 (20% of table size). This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_analyze_scale_factor (floating point)`

Specifies a fraction of the table size to add to `autovacuum_analyze_threshold` when deciding whether to trigger an `ANALYZE`. The default is 0.1 (10% of table size). This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_freeze_max_age (integer 64)`

Specifies the maximum age (in transactions) that a table's `pg_class.relFrozenxid` field can attain before a `VACUUM` operation is forced to shrink `pg_xact`. Note that the system will launch autovacuum processes for this purpose even when autovacuum is otherwise disabled. The default is 10 billion.

Vacuum also allows removal of old files from the `pg_xact` subdirectory. This parameter can only be set at server start, but the setting can be reduced for individual tables by changing table storage parameters. For more information see [Section 24.1.5](#).

`autovacuum_multixact_freeze_max_age (integer 64)`

Specifies the maximum age (in multixacts) that a table's `pg_class.relminmxid` field can attain before a `VACUUM` operation is forced to shrink `pg_multixact`. Note that the system will launch autovacuum processes for this purpose even when autovacuum is otherwise disabled. The default is 20 billion.

This parameter can only be set at server start, but the setting can be reduced for individual tables by changing table storage parameters. For more information see [Section 24.1.5.1](#).

`autovacuum_vacuum_cost_delay (floating point)`

Specifies the cost delay value that will be used in automatic `VACUUM` operations. If -1 is specified, the regular `vacuum_cost_delay` value will be used. If this value is specified without units, it is taken as milliseconds. The default value is 2 milliseconds. This parameter can only be set in the `post-`

`gresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`autovacuum_vacuum_cost_limit` (integer)

Specifies the cost limit value that will be used in automatic `VACUUM` operations. If -1 is specified (which is the default), the regular `vacuum_cost_limit` value will be used. Note that the value is distributed proportionally among the running autovacuum workers, if there is more than one, so that the sum of the limits for each worker does not exceed the value of this variable. This parameter can only be set in the `postgresql.conf` file or on the server command line; but the setting can be overridden for individual tables by changing table storage parameters.

`max_parallel_autovacuum_workers` (integer)

Sets the maximum number of additional autovacuum worker processes per cluster that can be used for parallel table vacuuming. Must be less than `autovacuum_max_workers`. The default is 0, which means no parallel table vacuuming.

19.11. Background Freezer

These settings control the behavior of background freezer. Refer to [Section 24.2](#) for more information.

`enable_background_freezer` (boolean)

Enables or disables sending messages to the background freezer worker using a UDP socket. This configuration parameter can be set either in the `postgresql.conf` file or in a session by superusers. The default value is `off`.

`bghfreezer_pages_per_rel_to_clean` (integer)

Sets the number of pages that the background freezer worker accumulates in a queue for each table to start vacuuming. For example, if `bghfreezer_pages_per_rel_to_clean=128`, the worker accumulates at least 128 pages in a queue and starts processing them after the transaction is completed. When set to 0, queued pages start to be processed at the end of each transaction. This configuration parameter can be set either in the `postgresql.conf` file or by sending the `SIGHUP` signal to the main server process. The default value is 128.

`bghfreezer_max_pages_per_rel` (integer)

Sets the maximum number of pages that can be queued for vacuuming per table. When this number is reached, page accumulation stops and new pages evict old pages from the queue. This configuration parameter can be set either in the `postgresql.conf` file or by sending the `SIGHUP` signal to the main server process. The default value is 32768.

19.12. Client Connection Defaults

19.12.1. Statement Behavior

`client_min_messages` (enum)

Controls which [message levels](#) are sent to the client. Valid values are `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `LOG`, `NOTICE`, `WARNING`, and `ERROR`. Each level includes all the levels that follow it. The later the level, the fewer messages are sent. The default is `NOTICE`. Note that `LOG` has a different rank here than in [log_min_messages](#).

`INFO` level messages are always sent to the client.

`search_path` (string)

This variable specifies the order in which schemas are searched when an object (table, data type, function, etc.) is referenced by a simple name with no schema specified. When there are objects of identical names in different schemas, the one found first in the search path is used. An object that

is not in any of the schemas in the search path can only be referenced by specifying its containing schema with a qualified (dotted) name.

The value for `search_path` must be a comma-separated list of schema names. Any name that is not an existing schema, or is a schema for which the user does not have `USAGE` permission, is silently ignored.

If one of the list items is the special name `$user`, then the schema having the name returned by `CURRENT_USER` is substituted, if there is such a schema and the user has `USAGE` permission for it. (If not, `$user` is ignored.)

The system catalog schema, `pg_catalog`, is always searched, whether it is mentioned in the path or not. If it is mentioned in the path then it will be searched in the specified order. If `pg_catalog` is not in the path then it will be searched *before* searching any of the path items.

Likewise, the current session's temporary-table schema, `pg_temp_nnn`, is always searched if it exists. It can be explicitly listed in the path by using the alias `pg_temp`. If it is not listed in the path then it is searched first (even before `pg_catalog`). However, the temporary schema is only searched for relation (table, view, sequence, etc.) and data type names. It is never searched for function or operator names.

When objects are created without specifying a particular target schema, they will be placed in the first valid schema named in `search_path`. An error is reported if the search path is empty.

The default value for this parameter is `"$user", public`. This setting supports shared use of a database (where no users have private schemas, and all share use of `public`), private per-user schemas, and combinations of these. Other effects can be obtained by altering the default search path setting, either globally or per-user.

For more information on schema handling, see [Section 5.10](#). In particular, the default configuration is suitable only when the database has a single user or a few mutually-trusting users.

The current effective value of the search path can be examined via the SQL function `current_schemas` (see [Section 9.27](#)). This is not quite the same as examining the value of `search_path`, since `current_schemas` shows how the items appearing in `search_path` were resolved.

`row_security` (boolean)

This variable controls whether to raise an error in lieu of applying a row security policy. When set to `on`, policies apply normally. When set to `off`, queries fail which would otherwise apply at least one policy. The default is `on`. Change to `off` where limited row visibility could cause incorrect results; for example, `pg_dump` makes that change by default. This variable has no effect on roles which bypass every row security policy, to wit, superusers and roles with the `BYPASSRLS` attribute.

For more information on row security policies, see [CREATE POLICY](#).

`default_table_access_method` (string)

This parameter specifies the default table access method to use when creating tables or materialized views if the `CREATE` command does not explicitly specify an access method, or when `SELECT ... INTO` is used, which does not allow specifying a table access method. The default is `heap`.

`default_tablespace` (string)

This variable specifies the default tablespace in which to create objects (tables and indexes) when a `CREATE` command does not explicitly specify a tablespace.

The value is either the name of a tablespace, or an empty string to specify using the default tablespace of the current database. If the value does not match the name of any existing tablespace, Postgres Pro will automatically use the default tablespace of the current database. If a nondefault tablespace is specified, the user must have `CREATE` privilege for it, or creation attempts will fail.

This variable is not used for temporary tables; for them, [temp_tablespaces](#) is consulted instead.

This variable is also not used when creating databases. By default, a new database inherits its tablespace setting from the template database it is copied from.

If this parameter is set to a value other than the empty string when a partitioned table is created, the partitioned table's tablespace will be set to that value, which will be used as the default tablespace for partitions created in the future, even if `default_tablespace` has changed since then.

For more information on tablespaces, see [Section 22.7](#).

`default_toast_compression` (enum)

This variable sets the default [TOAST](#) compression method for values of compressible columns. (This can be overridden for individual columns by setting the `COMPRESSION` column option in `CREATE TABLE` or `ALTER TABLE`.) The supported compression methods are `pglz` and (if Postgres Pro was compiled with `--with-lz4`) `lz4`. The default is `lz4`.

`temp_tablespaces` (string)

This variable specifies tablespaces in which to create temporary objects (temp tables and indexes on temp tables) when a `CREATE` command does not explicitly specify a tablespace. Temporary files for purposes such as sorting large data sets are also created in these tablespaces.

The value is a list of names of tablespaces. When there is more than one name in the list, Postgres Pro chooses a random member of the list each time a temporary object is to be created; except that within a transaction, successively created temporary objects are placed in successive tablespaces from the list. If the selected element of the list is an empty string, Postgres Pro will automatically use the default tablespace of the current database instead.

When `temp_tablespaces` is set interactively, specifying a nonexistent tablespace is an error, as is specifying a tablespace for which the user does not have `CREATE` privilege. However, when using a previously set value, nonexistent tablespaces are ignored, as are tablespaces for which the user lacks `CREATE` privilege. In particular, this rule applies when using a value set in `postgresql.conf`.

The default value is an empty string, which results in all temporary objects being created in the default tablespace of the current database.

See also [default_tablespace](#).

`check_function_bodies` (boolean)

This parameter is normally on. When set to `off`, it disables validation of the routine body string during [CREATE FUNCTION](#) and [CREATE PROCEDURE](#). Disabling validation avoids side effects of the validation process, in particular preventing false positives due to problems such as forward references. Set this parameter to `off` before loading functions on behalf of other users; `pg_dump` does so automatically.

`default_transaction_isolation` (enum)

Each SQL transaction has an isolation level, which can be either “read uncommitted”, “read committed”, “repeatable read”, or “serializable”. This parameter controls the default isolation level of each new transaction. The default is “read committed”.

Consult [Chapter 13](#) and [SET TRANSACTION](#) for more information.

`default_transaction_read_only` (boolean)

A read-only SQL transaction cannot alter non-temporary tables. This parameter controls the default read-only status of each new transaction. The default is `off` (read/write).

Consult [SET TRANSACTION](#) for more information.

`default_transaction_deferrable` (boolean)

When running at the `serializable` isolation level, a deferrable read-only SQL transaction may be delayed before it is allowed to proceed. However, once it begins executing it does not incur any of

the overhead required to ensure serializability; so serialization code will have no reason to force it to abort because of concurrent updates, making this option suitable for long-running read-only transactions.

This parameter controls the default deferrable status of each new transaction. It currently has no effect on read-write transactions or those operating at isolation levels lower than `serializable`. The default is `off`.

Consult [SET TRANSACTION](#) for more information.

`transaction_isolation` (enum)

This parameter reflects the current transaction's isolation level. At the beginning of each transaction, it is set to the current value of [default_transaction_isolation](#). Any subsequent attempt to change it is equivalent to a [SET TRANSACTION](#) command.

`transaction_read_only` (boolean)

This parameter reflects the current transaction's read-only status. At the beginning of each transaction, it is set to the current value of [default_transaction_read_only](#). Any subsequent attempt to change it is equivalent to a [SET TRANSACTION](#) command.

`transaction_deferrable` (boolean)

This parameter reflects the current transaction's deferrability status. At the beginning of each transaction, it is set to the current value of [default_transaction_deferrable](#). Any subsequent attempt to change it is equivalent to a [SET TRANSACTION](#) command.

`session_replication_role` (enum)

Controls firing of replication-related triggers and rules for the current session. Possible values are `origin` (the default), `replica` and `local`. Setting this parameter results in discarding any previously cached query plans. Only superusers and users with the appropriate `SET` privilege can change this setting.

The intended use of this setting is that logical replication systems set it to `replica` when they are applying replicated changes. The effect of that will be that triggers and rules (that have not been altered from their default configuration) will not fire on the replica. See the [ALTER TABLE](#) clauses `ENABLE TRIGGER` and `ENABLE RULE` for more information.

Postgres Pro treats the settings `origin` and `local` the same internally. Third-party replication systems may use these two values for their internal purposes, for example using `local` to designate a session whose changes should not be replicated.

Since foreign keys are implemented as triggers, setting this parameter to `replica` also disables all foreign key checks, which can leave data in an inconsistent state if improperly used.

`statement_timeout` (integer)

Abort any statement that takes more than the specified amount of time. If `log_min_error_statement` is set to `ERROR` or lower, the statement that timed out will also be logged. If this value is specified without units, it is taken as milliseconds. A value of zero (the default) disables the timeout.

The timeout is measured from the time a command arrives at the server until it is completed by the server. If multiple SQL statements appear in a single simple-query message, the timeout is applied to each statement separately. (Postgres Pro versions before 13 usually treated the timeout as applying to the whole query string.) In extended query protocol, the timeout starts running when any query-related message (Parse, Bind, Execute, Describe) arrives, and it is canceled by completion of an Execute or Sync message.

Setting `statement_timeout` in `postgresql.conf` is not recommended because it would affect all sessions.

`transaction_timeout` (integer)

Terminate any session that spans longer than the specified amount of time in a transaction. The limit applies both to explicit transactions (started with `BEGIN`) and to an implicitly started transaction corresponding to a single statement. If this value is specified without units, it is taken as milliseconds. A value of zero (the default) disables the timeout.

If `transaction_timeout` is shorter or equal to `idle_in_transaction_session_timeout` or `statement_timeout` then the longer timeout is ignored.

Setting `transaction_timeout` in `postgresql.conf` is not recommended because it would affect all sessions.

Note

Prepared transactions are not subject to this timeout.

`lock_timeout` (integer)

Abort any statement that waits longer than the specified amount of time while attempting to acquire a lock on a table, index, row, or other database object. The time limit applies separately to each lock acquisition attempt. The limit applies both to explicit locking requests (such as `LOCK TABLE`, or `SELECT FOR UPDATE` without `NOWAIT`) and to implicitly-acquired locks. If this value is specified without units, it is taken as milliseconds. A value of zero (the default) disables the timeout.

Unlike `statement_timeout`, this timeout can only occur while waiting for locks. Note that if `statement_timeout` is nonzero, it is rather pointless to set `lock_timeout` to the same or larger value, since the statement timeout would always trigger first. If `log_min_error_statement` is set to `ERROR` or lower, the statement that timed out will be logged.

Setting `lock_timeout` in `postgresql.conf` is not recommended because it would affect all sessions.

`idle_in_transaction_session_timeout` (integer)

Terminate any session that has been idle (that is, waiting for a client query) within an open transaction for longer than the specified amount of time. If this value is specified without units, it is taken as milliseconds. A value of zero (the default) disables the timeout.

This option can be used to ensure that idle sessions do not hold locks for an unreasonable amount of time. Even when no significant locks are held, an open transaction prevents vacuuming away recently-dead tuples that may be visible only to this transaction; so remaining idle for a long time can contribute to table bloat. See [Section 24.1](#) for more details.

`idle_session_timeout` (integer)

Terminate any session that has been idle (that is, waiting for a client query), but not within an open transaction, for longer than the specified amount of time. If this value is specified without units, it is taken as milliseconds. A value of zero (the default) disables the timeout.

Unlike the case with an open transaction, an idle session without a transaction imposes no large costs on the server, so there is less need to enable this timeout than `idle_in_transaction_session_timeout`.

Be wary of enforcing this timeout on connections made through connection-pooling software or other middleware, as such a layer may not react well to unexpected connection closure. It may be helpful to enable this timeout only for interactive sessions, perhaps by applying it only to particular users.

`vacuum_freeze_table_age` (integer 64)

`VACUUM` performs an aggressive scan if the table's `pg_class.relfrozensid` field has reached the age specified by this setting. An aggressive scan differs from a regular `VACUUM` in that it visits every

page that might contain unfrozen XIDs or MXIDs, not just those that might contain dead tuples. The default is 150 million transactions. Although users can set this value anywhere from zero to $2^{63} - 1$, `VACUUM` will silently limit the effective value to 95% of [autovacuum_freeze_max_age](#), so that a periodic manual `VACUUM` has a chance to run before an autovacuum to shrink `pg_xact` and `pg_multixact` is launched for the table. This value is also used as a “soft” limit for advancing `datfrozenxid`, helping to prevent `pg_multixact` bloat without immediately triggering autovacuum, which would require an exclusive database lock. For more information see [Section 24.1.5](#).

`vacuum_freeze_min_age` (integer 64)

Specifies the cutoff age (in transactions) that `VACUUM` should use to decide whether to trigger freezing of pages that have an older XID. The default is 50 million transactions. Although users can set this value anywhere from zero to $2^{63} - 1$, `VACUUM` will silently limit the effective value to half the value of [autovacuum_freeze_max_age](#), so that there is not an unreasonably short time between forced autovacums. For more information see [Section 24.1.5](#).

`vacuum_failsafe_age` (integer)

Specifies the maximum age (in transactions) that a table's `pg_class.relFrozenxid` field can attain before `VACUUM` takes extraordinary measures to avoid system-wide transaction ID wraparound failure. This is `VACUUM`'s strategy of last resort. The failsafe typically triggers when an autovacuum to prevent transaction ID wraparound has already been running for some time, though it's possible for the failsafe to trigger during any `VACUUM`.

When the failsafe is triggered, any cost-based delay that is in effect will no longer be applied, further non-essential maintenance tasks (such as index vacuuming) are bypassed, and any [Buffer Access Strategy](#) in use will be disabled resulting in `VACUUM` being free to make use of all of [shared buffers](#).

The default is 1.6 billion transactions. Although users can set this value anywhere from zero to 2.1 billion, `VACUUM` will silently adjust the effective value to no less than 105% of [autovacuum_freeze_max_age](#).

`vacuum_multixact_freeze_table_age` (integer 64)

`VACUUM` performs an aggressive scan if the table's `pg_class.relminmxid` field has reached the age specified by this setting. An aggressive scan differs from a regular `VACUUM` in that it visits every page that might contain unfrozen XIDs or MXIDs, not just those that might contain dead tuples. The default is 150 million multixacts. Although users can set this value anywhere from zero to $2^{63} - 1$, `VACUUM` will silently limit the effective value to 95% of [autovacuum_multixact_freeze_max_age](#), so that a periodic manual `VACUUM` has a chance to run before an autovacuum to shrink `pg_xact` and `pg_multixact` is launched for the table. This value is also used as a “soft” limit for advancing `datminmxid`, helping to prevent `pg_multixact` bloat without immediately triggering autovacuum, which would require an exclusive database lock. For more information see [Section 24.1.5](#).

`vacuum_multixact_freeze_min_age` (integer 64)

Specifies the cutoff age (in multixacts) that `VACUUM` should use to decide whether to trigger freezing of pages with an older multixact ID. The default is 5 million multixacts. Although users can set this value anywhere from zero to $2^{63} - 1$, `VACUUM` will silently limit the effective value to half the value of [autovacuum_multixact_freeze_max_age](#), so that there is not an unreasonably short time between forced autovacums. For more information see [Section 24.1.5.1](#).

`vacuum_multixact_failsafe_age` (integer)

Specifies the maximum age (in multixacts) that a table's `pg_class.relminmxid` field can attain before `VACUUM` takes extraordinary measures to avoid system-wide multixact ID wraparound failure. This is `VACUUM`'s strategy of last resort. The failsafe typically triggers when an autovacuum to prevent transaction ID wraparound has already been running for some time, though it's possible for the failsafe to trigger during any `VACUUM`.

When the failsafe is triggered, any cost-based delay that is in effect will no longer be applied, and further non-essential maintenance tasks (such as index vacuuming) are bypassed.

The default is 1.6 billion multixacts. Although users can set this value anywhere from zero to 2.1 billion, `VACUUM` will silently adjust the effective value to no less than 105% of `autovacuum_multixact_freeze_max_age`.

`vacuum_update_datfrozenxid_only_when_needed` (boolean)

`datfrozenxid` is usually updated for tables during or after `VACUUM` in a sequential scan of the `pg_class` table. The update process lasts long when many tables are vacuumed. When `vacuum_update_datfrozenxid_only_when_needed` is set to on, `datfrozenxid` for a table only gets updated in a sequential scan if it was actually updated. This can significantly improve the performance when vacuuming a large number of tables. The default is `off`.

`bytea_output` (enum)

Sets the output format for values of type `bytea`. Valid values are `hex` (the default) and `escape` (the traditional Postgres Pro format). See [Section 8.4](#) for more information. The `bytea` type always accepts both formats on input, regardless of this setting.

`xmlbinary` (enum)

Sets how binary values are to be encoded in XML. This applies for example when `bytea` values are converted to XML by the functions `xmlelement` or `xmlforest`. Possible values are `base64` and `hex`, which are both defined in the XML Schema standard. The default is `base64`. For further information about XML-related functions, see [Section 9.15](#).

The actual choice here is mostly a matter of taste, constrained only by possible restrictions in client applications. Both methods support all possible values, although the hex encoding will be somewhat larger than the base64 encoding.

`xmloption` (enum)

Sets whether `DOCUMENT` or `CONTENT` is implicit when converting between XML and character string values. See [Section 8.13](#) for a description of this. Valid values are `DOCUMENT` and `CONTENT`. The default is `CONTENT`.

According to the SQL standard, the command to set this option is

```
SET XML OPTION { DOCUMENT | CONTENT };
```

This syntax is also available in Postgres Pro.

`gin_pending_list_limit` (integer)

Sets the maximum size of a GIN index's pending list, which is used when `fastupdate` is enabled. If the list grows larger than this maximum size, it is cleaned up by moving the entries in it to the index's main GIN data structure in bulk. If this value is specified without units, it is taken as kilobytes. The default is four megabytes (4MB). This setting can be overridden for individual GIN indexes by changing index storage parameters. See [Section 65.4.4.1](#) and [Section 65.4.5](#) for more information.

`creatorole_self_grant` (string)

If a user who has `CREATEROLE` but not `SUPERUSER` creates a role, and if this is set to a non-empty value, the newly-created role will be granted to the creating user with the options specified. The value must be `set`, `inherit`, or a comma-separated list of these. The default value is an empty string, which disables the feature.

The purpose of this option is to allow a `CREATEROLE` user who is not a superuser to automatically inherit, or automatically gain the ability to `SET ROLE` to, any created users. Since a `CREATEROLE` user is always implicitly granted `ADMIN OPTION` on created roles, that user could always execute a `GRANT` statement that would achieve the same effect as this setting. However, it can be convenient for usability reasons if the grant happens automatically. A superuser automatically inherits the privileges of every role and can always `SET ROLE` to any role, and this setting can be used to produce a similar behavior for `CREATEROLE` users for users which they create.

`event_triggers` (boolean)

Allow temporarily disabling execution of event triggers in order to troubleshoot and repair faulty event triggers. All event triggers will be disabled by setting it to `false`. Setting the value to `true` allows all event triggers to fire, this is the default value. Only superusers and users with the appropriate `SET` privilege can change this setting.

`restrict_nonsystem_relation_kind` (string)

Set relation kinds for which access to non-system relations is prohibited. The value takes the form of a comma-separated list of relation kinds. Currently, the supported relation kinds are `view` and `foreign-table`.

`skip_temp_rel_lock` (boolean)

Controls locking behavior for temporary relations. When set to `on`, locking is skipped for temporary relations. This could improve performance for applications that frequently use such relations. The default is `off`.

19.12.2. Locale and Formatting

`DateStyle` (string)

Sets the display format for date and time values, as well as the rules for interpreting ambiguous date input values. For historical reasons, this variable contains two independent components: the output format specification (`ISO`, `Postgres`, `SQL`, or `German`) and the input/output specification for year/month/day ordering (`DMY`, `MDY`, or `YMD`). These can be set separately or together. The keywords `Euro` and `European` are synonyms for `DMY`; the keywords `US`, `NonEuro`, and `NonEuropean` are synonyms for `MDY`. See [Section 8.5](#) for more information. The built-in default is `ISO`, `MDY`, but `initdb` will initialize the configuration file with a setting that corresponds to the behavior of the chosen `lc_time` locale.

`IntervalStyle` (enum)

Sets the display format for interval values. The value `sql_standard` will produce output matching SQL standard interval literals. The value `postgres` (which is the default) will produce output matching PostgreSQL releases prior to 8.4 when the `DateStyle` parameter was set to `ISO`. The value `postgres_verbose` will produce output matching PostgreSQL releases prior to 8.4 when the `DateStyle` parameter was set to non-ISO output. The value `iso_8601` will produce output matching the time interval “format with designators” defined in section 4.4.3.2 of ISO 8601.

The `IntervalStyle` parameter also affects the interpretation of ambiguous interval input. See [Section 8.5.4](#) for more information.

`TimeZone` (string)

Sets the time zone for displaying and interpreting time stamps. The built-in default is `GMT`, but that is typically overridden in `postgresql.conf`; `initdb` will install a setting there corresponding to its system environment. See [Section 8.5.3](#) for more information.

`timezone_abbreviations` (string)

Sets the collection of time zone abbreviations that will be accepted by the server for datetime input. The default is `'Default'`, which is a collection that works in most of the world; there are also `'Australia'` and `'India'`, and other collections can be defined for a particular installation. See [Section B.4](#) for more information.

`extra_float_digits` (integer)

This parameter adjusts the number of digits used for textual output of floating-point values, including `float4`, `float8`, and geometric data types.

If the value is 1 (the default) or above, float values are output in shortest-precise format; see [Section 8.1.3](#). The actual number of digits generated depends only on the value being output, not on the

value of this parameter. At most 17 digits are required for `float8` values, and 9 for `float4` values. This format is both fast and precise, preserving the original binary float value exactly when correctly read. For historical compatibility, values up to 3 are permitted.

If the value is zero or negative, then the output is rounded to a given decimal precision. The precision used is the standard number of digits for the type (`FLT_DIG` or `DBL_DIG` as appropriate) reduced according to the value of this parameter. (For example, specifying `-1` will cause `float4` values to be output rounded to 5 significant digits, and `float8` values rounded to 14 digits.) This format is slower and does not preserve all the bits of the binary float value, but may be more human-readable.

Note

The meaning of this parameter, and its default value, changed in Postgres Pro 12; see [Section 8.1.3](#) for further discussion.

`client_encoding (string)`

Sets the client-side encoding (character set). The default is to use the database encoding. The character sets supported by the Postgres Pro server are described in [Section 23.3.1](#).

`lc_messages (string)`

Sets the language in which messages are displayed. Acceptable values are system-dependent; see [Section 23.1](#) for more information. If this variable is set to the empty string (which is the default) then the value is inherited from the execution environment of the server in a system-dependent way.

On some systems, this locale category does not exist. Setting this variable will still work, but there will be no effect. Also, there is a chance that no translated messages for the desired language exist. In that case you will continue to see the English messages.

Only superusers and users with the appropriate `SET` privilege can change this setting.

`lc_monetary (string)`

Sets the locale to use for formatting monetary amounts, for example with the `to_char` family of functions. Acceptable values are system-dependent; see [Section 23.1](#) for more information. If this variable is set to the empty string (which is the default) then the value is inherited from the execution environment of the server in a system-dependent way.

`lc_numeric (string)`

Sets the locale to use for formatting numbers, for example with the `to_char` family of functions. Acceptable values are system-dependent; see [Section 23.1](#) for more information. If this variable is set to the empty string (which is the default) then the value is inherited from the execution environment of the server in a system-dependent way.

`lc_time (string)`

Sets the locale to use for formatting dates and times, for example with the `to_char` family of functions. Acceptable values are system-dependent; see [Section 23.1](#) for more information. If this variable is set to the empty string (which is the default) then the value is inherited from the execution environment of the server in a system-dependent way.

`icu_validation_level (enum)`

When ICU locale validation problems are encountered, controls which [message level](#) is used to report the problem. Valid values are `DISABLED`, `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `INFO`, `NOTICE`, `WARNING`, `ERROR`, and `LOG`.

If set to `DISABLED`, does not report validation problems at all. Otherwise reports problems at the given message level. The default is `WARNING`.

`default_text_search_config (string)`

Selects the text search configuration that is used by those variants of the text search functions that do not have an explicit argument specifying the configuration. See [Chapter 12](#) for further information. The built-in default is `pg_catalog.simple`, but `initdb` will initialize the configuration file with a setting that corresponds to the chosen `lc_ctype` locale, if a configuration matching that locale can be identified.

19.12.3. Shared Library Preloading

Several settings are available for preloading shared libraries into the server, in order to load additional functionality or achieve performance benefits. For example, a setting of `'$libdir/mylib'` would cause `mylib.so` (or on some platforms, `mylib.sl`) to be preloaded from the installation's standard library directory. The differences between the settings are when they take effect and what privileges are required to change them.

Postgres Pro procedural language libraries can be preloaded in this way, typically by using the syntax `'$libdir/plXXX'` where `XXX` is `pgsql`, `perl`, `tcl`, or `python`.

Only shared libraries specifically intended to be used with Postgres Pro can be loaded this way. Every Postgres Pro-supported library has a “magic block” that is checked to guarantee compatibility. For this reason, non-Postgres Pro libraries cannot be loaded in this way. You might be able to use operating-system facilities such as `LD_PRELOAD` for that.

In general, refer to the documentation of a specific module for the recommended way to load that module.

`local_preload_libraries (string)`

This variable specifies one or more shared libraries that are to be preloaded at connection start. It contains a comma-separated list of library names, where each name is interpreted as for the [LOAD](#) command. Whitespace between entries is ignored; surround a library name with double quotes if you need to include whitespace or commas in the name. The parameter value only takes effect at the start of the connection. Subsequent changes have no effect. If a specified library is not found, the connection attempt will fail.

This option can be set by any user. Because of that, the libraries that can be loaded are restricted to those appearing in the `plugins` subdirectory of the installation's standard library directory. (It is the database administrator's responsibility to ensure that only “safe” libraries are installed there.) Entries in `local_preload_libraries` can specify this directory explicitly, for example `$libdir/plugins/mylib`, or just specify the library name — `mylib` would have the same effect as `$libdir/plugins/mylib`.

The intent of this feature is to allow unprivileged users to load debugging or performance-measurement libraries into specific sessions without requiring an explicit `LOAD` command. To that end, it would be typical to set this parameter using the `PGOPTIONS` environment variable on the client or by using `ALTER ROLE SET`.

However, unless a module is specifically designed to be used in this way by non-superusers, this is usually not the right setting to use. Look at [session_preload_libraries](#) instead.

`session_preload_libraries (string)`

This variable specifies one or more shared libraries that are to be preloaded at connection start. It contains a comma-separated list of library names, where each name is interpreted as for the [LOAD](#) command. Whitespace between entries is ignored; surround a library name with double quotes if you need to include whitespace or commas in the name. The parameter value only takes effect at the start of the connection. Subsequent changes have no effect. If a specified library is not found, the connection attempt will fail. Only superusers and users with the appropriate `SET` privilege can change this setting.

The intent of this feature is to allow debugging or performance-measurement libraries to be loaded into specific sessions without an explicit `LOAD` command being given. For example, [auto_explain](#) could

be enabled for all sessions under a given user name by setting this parameter with `ALTER ROLE SET`. Also, this parameter can be changed without restarting the server (but changes only take effect when a new session is started), so it is easier to add new modules this way, even if they should apply to all sessions.

Unlike [shared_preload_libraries](#), there is no large performance advantage to loading a library at session start rather than when it is first used. There is some advantage, however, when connection pooling is used.

`shared_preload_libraries (string)`

This variable specifies one or more shared libraries to be preloaded at server start. It contains a comma-separated list of library names, where each name is interpreted as for the `LOAD` command. Whitespace between entries is ignored; surround a library name with double quotes if you need to include whitespace or commas in the name. This parameter can only be set at server start. If a specified library is not found, the server will fail to start.

For distributed system specifics, see [shared_preload_libraries](#).

Some libraries need to perform certain operations that can only take place at postmaster start, such as allocating shared memory, reserving light-weight locks, or starting background workers. Those libraries must be loaded at server start through this parameter. See the documentation of each library for details.

Other libraries can also be preloaded. By preloading a shared library, the library startup time is avoided when the library is first used. However, the time to start each new server process might increase slightly, even if that process never uses the library. So this parameter is recommended only for libraries that will be used in most sessions. Also, changing this parameter requires a server restart, so this is not the right setting to use for short-term debugging tasks, say. Use [session_preload_libraries](#) for that instead.

Note

On Windows hosts, preloading a library at server start will not reduce the time required to start each new server process; each server process will re-load all preload libraries. However, `shared_preload_libraries` is still useful on Windows hosts for libraries that need to perform operations at postmaster start time.

`jit_provider (string)`

This variable is the name of the JIT provider library to be used (see [Section 30.4.2](#)). The default is `llvmjit`. This parameter can only be set at server start.

If set to a non-existent library, JIT will not be available, but no error will be raised. This allows JIT support to be installed separately from the main Postgres Pro package.

19.12.4. Other Defaults

`dynamic_library_path (string)`

If a dynamically loadable module needs to be opened and the file name specified in the `CREATE FUNCTION` or `LOAD` command does not have a directory component (i.e., the name does not contain a slash), the system will search this path for the required file.

The value for `dynamic_library_path` must be a list of absolute directory paths separated by colons (or semi-colons on Windows). If a list element starts with the special string `$libdir`, the compiled-in Postgres Pro package library directory is substituted for `$libdir`; this is where the modules provided by the standard Postgres Pro distribution are installed. (Use `pg_config --pkglibdir` to find out the name of this directory.) For example:

```
dynamic_library_path = '/usr/local/lib/postgresql:/home/my_project/lib:$libdir'
```

or, in a Windows environment:

```
dynamic_library_path = 'C:\tools\postgresql;H:\my_project\lib;$libdir'
```

The default value for this parameter is '\$libdir'. If the value is set to an empty string, the automatic path search is turned off.

This parameter can be changed at run time by superusers and users with the appropriate `SET` privilege, but a setting done that way will only persist until the end of the client connection, so this method should be reserved for development purposes. The recommended way to set this parameter is in the `postgresql.conf` configuration file.

`gin_fuzzy_search_limit` (integer)

Soft upper limit of the size of the set returned by GIN index scans. For more information see [Section 65.4.5](#).

19.13. Lock Management

`deadlock_timeout` (integer)

This is the amount of time to wait on a lock before checking to see if there is a deadlock condition. The check for deadlock is relatively expensive, so the server doesn't run it every time it waits for a lock. We optimistically assume that deadlocks are not common in production applications and just wait on the lock for a while before checking for a deadlock. Increasing this value reduces the amount of time wasted in needless deadlock checks, but slows down reporting of real deadlock errors. If this value is specified without units, it is taken as milliseconds. The default is one second (1s), which is probably about the smallest value you would want in practice. On a heavily loaded server you might want to raise it. Ideally the setting should exceed your typical transaction time, so as to improve the odds that a lock will be released before the waiter decides to check for deadlock. Only superusers and users with the appropriate `SET` privilege can change this setting.

When [log_lock_waits](#) is set, this parameter also determines the amount of time to wait before a log message is issued about the lock wait. If you are trying to investigate locking delays you might want to set a shorter than normal `deadlock_timeout`.

`max_locks_per_transaction` (integer)

The shared lock table has space for `max_locks_per_transaction` objects (e.g., tables) per server process or prepared transaction; hence, no more than this many distinct objects can be locked at any one time. This parameter limits the average number of object locks used by each transaction; individual transactions can lock more objects as long as the locks of all transactions fit in the lock table. This is *not* the number of rows that can be locked; that value is unlimited. The default, 64, has historically proven sufficient, but you might need to raise this value if you have queries that touch many different tables in a single transaction, e.g., query of a parent table with many children. This parameter can only be set at server start.

When running a standby server, you must set this parameter to have the same or higher value as on the primary server. Otherwise, queries will not be allowed in the standby server.

`max_pred_locks_per_transaction` (integer)

The shared predicate lock table has space for `max_pred_locks_per_transaction` objects (e.g., tables) per server process or prepared transaction; hence, no more than this many distinct objects can be locked at any one time. This parameter limits the average number of object locks used by each transaction; individual transactions can lock more objects as long as the locks of all transactions fit in the lock table. This is *not* the number of rows that can be locked; that value is unlimited. The default, 64, has historically proven sufficient, but you might need to raise this value if you have clients that touch many different tables in a single serializable transaction. This parameter can only be set at server start.

`max_pred_locks_per_relation` (integer)

This controls how many pages or tuples of a single relation can be predicate-locked before the lock is promoted to covering the whole relation. Values greater than or equal to zero mean an absolute limit, while negative values mean `max_pred_locks_per_transaction` divided by the absolute value of this setting. The default is -2, which keeps the behavior from previous versions of Postgres Pro. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`max_pred_locks_per_page` (integer)

This controls how many rows on a single page can be predicate-locked before the lock is promoted to covering the whole page. The default is 2. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`lwlock_shared_limit` (integer)

Specifies the number of sequential shared LWLocks held before switching to the *fair mode*. In the fair mode, would-be shared lockers are added to the queue instead of acquiring the lock immediately, which allows to acquire an exclusive LWLock within reasonable time and thus avoid lock starvation. Setting this option to small values allows to acquire exclusive locks much faster, but can cause significant overhead. The larger the value, the lower the overhead and the speedup. As a rule of thumb, setting this option to 16 provides good speedup with minimal overhead. The default is 0 (the option is disabled). This parameter can only be set in the `postgresql.conf` file or on the server command line.

`log2_num_lock_partitions` (integer)

This controls how many partitions the shared lock tables are divided into. Number of partitions is calculated by raising 2 to the power of this parameter. The default value is 4, which corresponds to 16 partitions, and the maximum is 8. This parameter can only be set in the `postgresql.conf` file or on the server command line.

19.13.1. CSN Based Snapshot

By default, The snapshots in PostgreSQL uses the XID (TransactionID) to identify the status of the transaction, the in-progress transactions, and the future transactions for all its visibility calculations.

PostgreSQL also provides the CSN (commit-sequence-number) based mechanism to identify the past-transactions and the ones that are yet to be started/committed.

`enable_csn_snapshot` (boolean)

Enable/disable the CSN based transaction visibility tracking for the snapshot.

For distributed system specifics, see `enable_csn_snapshot`.

PostgreSQL uses the clock timestamp as a CSN, so enabling the CSN based snapshots can be useful for implementing the global snapshots and global transaction visibility.

when enabled PostgreSQL creates `pg_csn` directory under `PGDATA` to keep the track of CSN and XID mappings.

The default value is off.

19.14. Version and Platform Compatibility

19.14.1. Previous Postgres Pro Versions

`array_nulls` (boolean)

This controls whether the array input parser recognizes unquoted `NULL` as specifying a null array element. By default, this is `on`, allowing array values containing null values to be entered. However, PostgreSQL versions before 8.2 did not support null values in arrays, and therefore would treat `NULL` as specifying a normal array element with the string value "NULL". For backward compatibility with applications that require the old behavior, this variable can be turned `off`.

Note that it is possible to create array values containing null values even when this variable is `off`.

`backslash_quote` (enum)

This controls whether a quote mark can be represented by `\'` in a string literal. The preferred, SQL-standard way to represent a quote mark is by doubling it (`' '`) but Postgres Pro has historically also accepted `\'`. However, use of `\'` creates security risks because in some client character set encodings, there are multibyte characters in which the last byte is numerically equivalent to ASCII `\`. If client-side code does escaping incorrectly then an SQL-injection attack is possible. This risk can be prevented by making the server reject queries in which a quote mark appears to be escaped by a backslash. The allowed values of `backslash_quote` are `on` (allow `\'` always), `off` (reject always), and `safe_encoding` (allow only if client encoding does not allow ASCII `\` within a multibyte character). `safe_encoding` is the default setting.

Note that in a standard-conforming string literal, `\` just means `\` anyway. This parameter only affects the handling of non-standard-conforming literals, including escape string syntax (`E' ... '`).

`escape_string_warning` (boolean)

When `on`, a warning is issued if a backslash (`\`) appears in an ordinary string literal (`' ... '` syntax) and `standard_conforming_strings` is `off`. The default is `on`.

Applications that wish to use backslash as escape should be modified to use escape string syntax (`E' ... '`), because the default behavior of ordinary strings is now to treat backslash as an ordinary character, per SQL standard. This variable can be enabled to help locate code that needs to be changed.

`lo_compat_privileges` (boolean)

In PostgreSQL releases prior to 9.0, large objects did not have access privileges and were, therefore, always readable and writable by all users. Setting this variable to `on` disables the new privilege checks, for compatibility with prior releases. The default is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

Setting this variable does not disable all security checks related to large objects — only those for which the default behavior has changed in PostgreSQL 9.0.

`quote_all_identifiers` (boolean)

When the database generates SQL, force all identifiers to be quoted, even if they are not (currently) keywords. This will affect the output of `EXPLAIN` as well as the results of functions like `pg_get_viewdef`. See also the `--quote-all-identifiers` option of [pg_dump](#) and [pg_dumpall](#).

`nul_byte_replacement_on_import` (string)

Replace NUL bytes `'\0'` with the specified decimal code of an ASCII character while loading data using the `COPY FROM` command. Such a replacement may be required when transferring data from another DBMS since Postgres Pro does not allow NUL bytes in data. The specified ASCII code must not coincide with the `QUOTE` and `DELIMITER` characters used by `COPY FROM` as it may cause unexpected results. The default value is `'\0'`, so no replacement occurs.

`enable_large_mem_buffers` (boolean)

When set to `true`, allows to copy, dump and restore large `bytea` values exceeding 0.5 GB (up to 1 GB) and records with several `text` values exceeding 1 GB in total (up to 2 GB). The default value is `false`. Only superusers can change this setting and only in the current session using the `SET` command.

`standard_conforming_strings` (boolean)

This controls whether ordinary string literals (`' ... '`) treat backslashes literally, as specified in the SQL standard. Beginning in PostgreSQL 9.1, the default is `on` (prior releases defaulted to `off`). Applications can check this parameter to determine how string literals will be processed. The presence of this parameter can also be taken as an indication that the escape string syntax (`E' ... '`) is sup-

ported. Escape string syntax ([Section 4.1.2.2](#)) should be used if an application desires backslashes to be treated as escape characters.

`synchronize_seqscans` (boolean)

This allows sequential scans of large tables to synchronize with each other, so that concurrent scans read the same block at about the same time and hence share the I/O workload. When this is enabled, a scan might start in the middle of the table and then “wrap around” the end to cover all rows, so as to synchronize with the activity of scans already in progress. This can result in unpredictable changes in the row ordering returned by queries that have no `ORDER BY` clause. Setting this parameter to `off` ensures the pre-8.3 behavior in which a sequential scan always starts from the beginning of the table. The default is `on`.

19.14.2. Platform and Client Compatibility

`unicode_nul_character_replacement_in_jsonb` (integer)

Replace `'\u0000'` with the specified unicode character when calling a function processing JSONB since `jsonb*` functions do not accept `'\u0000'`. The replacing character is specified with the numeric value from 0 to 65535. Anyone can set the value of this configuration parameter anytime. The default value is 0, so no replacement occurs.

`transform_null_equals` (boolean)

When `on`, expressions of the form `expr = NULL` (or `NULL = expr`) are treated as `expr IS NULL`, that is, they return true if `expr` evaluates to the null value, and false otherwise. The correct SQL-spec-compliant behavior of `expr = NULL` is to always return null (unknown). Therefore this parameter defaults to `off`.

However, filtered forms in Microsoft Access generate queries that appear to use `expr = NULL` to test for null values, so if you use that interface to access the database you might want to turn this option `on`. Since expressions of the form `expr = NULL` always return the null value (using the SQL standard interpretation), they are not very useful and do not appear often in normal applications so this option does little harm in practice. But new users are frequently confused about the semantics of expressions involving null values, so this option is `off` by default.

Note that this option only affects the exact form `= NULL`, not other comparison operators or other expressions that are computationally equivalent to some expression involving the equals operator (such as `IN`). Thus, this option is not a general fix for bad programming.

Refer to [Section 9.2](#) for related information.

`allow_alter_system` (boolean)

When `allow_alter_system` is set to `off`, an error is returned if the `ALTER SYSTEM` command is executed. This parameter can only be set in the `postgresql.conf` file or on the server command line. The default value is `on`.

Note that this setting must not be regarded as a security feature. It only disables the `ALTER SYSTEM` command. It does not prevent a superuser from changing the configuration using other SQL commands. A superuser has many ways of executing shell commands at the operating system level, and can therefore modify `postgresql.auto.conf` regardless of the value of this setting.

Turning this setting off is intended for environments where the configuration of Postgres Pro is managed by some external tool. In such environments, a well intentioned superuser might *mistakenly* use `ALTER SYSTEM` to change the configuration instead of using the external tool. This might result in unintended behavior, such as the external tool overwriting the change at some later point in time when it updates the configuration. Setting this parameter to `off` can help avoid such mistakes.

This parameter only controls the use of `ALTER SYSTEM`. The settings stored in `postgresql.auto.conf` take effect even if `allow_alter_system` is set to `off`.

19.15. Memory Purge

Postgres Pro Shardman provides memory purge configuration parameters that can automatically replace data with zero bytes before it is deleted. For details, see [Section 31.1](#).

`wipe_file_on_delete` (boolean)

Replace a file with zero bytes before it is deleted.

Default: `on`

`wipe_heaptuple_on_delete` (boolean)

Replace a row version with zero bytes before it is deleted.

Default: `on`

`wipe_memctx_on_free` (boolean)

Clear random-access memory before it is freed. This parameter only controls the memory allocated as part of a particular context.

Default: `on`

`wipe_mem_on_free` (boolean)

Clear random-access memory before it is freed. This parameter controls the memory that does not belong to any context.

Default: `on`

`wipe_xlog_on_free` (boolean)

Replace WAL segments with zero bytes before they are deleted or overwritten.

Default: `on`

19.16. Data Compression

`cfs_gc` (boolean)

Enables/disables background garbage collection of compressed pages. Default: `on`

`cfs_gc_lock_file` (char)

Path to the lock file to be used to ensure running only one garbage collection worker at a time for several Postgres Pro servers when `cfs_gc_workers` = 1. This parameter can only be set at server start.

`cfs_gc_workers` (integer)

Number of CFS background garbage collection workers. This parameter can only be set at server start. The number of workers cannot exceed 100. Default: 1

`cfs_level` (integer)

CFS compression level: 0 — no compression, 1 — maximal speed. Other possible values depend on the specific compression algorithm used. For example, 9 for `zlib` or 19 for `zstd`. Default: 1

`cfs_gc_threshold` (integer)

Minimum percent of garbage blocks in the file required to start garbage collection. If you would like to manually defragment a relation with a smaller percent of garbage, you can temporarily set this parameter to a smaller value for your current session. Default: 30

`cfs_gc_period` (integer)

Time interval between CFS garbage collection iterations, in milliseconds. Default: 5 seconds

`cfs_compress_temp_relations` (bool)

When set to `true`, compresses temporary tables. For large temporary tables, this setting can save disk space and speed up execution as less data needs to be read. This parameter can only be set at server start. Default: `false`

`cfs_gc_delay` (integer)

Time interval for which garbage collection is paused after each file defragmentation, in milliseconds. Default: 0

`cfs_gc_respond_time` (integer)

Time interval that CFS waits for the lock to be released from a file processed by garbage collector before writing a warning to the log, in seconds. Default: 3600

`cfs_gc_time_bloom_query` (integer)

Reports the number of times the Bloom filter was queried for CFS files. This parameter is read-only.

`cfs_gc_time_bloom_passed` (integer)

Reports the number of elements that passed through the Bloom filter. This parameter is read-only.

`cfs_gc_time_bloom_false_positive` (integer)

Reports the number of false positives returned by the Bloom filter. This parameter is read-only.

`cfs_log_verbose` (boolean)

Controls the CFS log level of defragmentation messages. If set to `false`, the level is `INFO`. Defragmentation messages are written to the log file when `log_min_messages` has a level from `DEBUG5` to `INFO`.

If set to `true`, the level is set to `LOG`. Defragmentation messages are written to the log file when `log_min_messages` has a level from `DEBUG5` to `LOG`. In this case, they all have `LOG` as prefix, regardless of their level. See [message levels](#) for more information.

19.17. Error Handling

`exit_on_error` (boolean)

If on, any error will terminate the current session. By default, this is set to off, so that only `FATAL` errors will terminate the session.

`restart_after_crash` (boolean)

When set to on, which is the default, Postgres Pro will automatically reinitialize after a backend crash. Leaving this value set to on is normally the best way to maximize the availability of the database. However, in some circumstances, such as when Postgres Pro is being invoked by clusterware, it may be useful to disable the restart so that the clusterware can gain control and take any actions it deems appropriate.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`data_sync_retry` (boolean)

When set to off, which is the default, Postgres Pro will raise a `PANIC`-level error on failure to flush modified data files to the file system. This causes the database server to crash. This parameter can only be set at server start.

On some operating systems, the status of data in the kernel's page cache is unknown after a write-back failure. In some cases it might have been entirely forgotten, making it unsafe to retry; the second attempt may be reported as successful, when in fact the data has been lost. In these circumstances, the only way to avoid data loss is to recover from the WAL after any failure is reported, preferably after investigating the root cause of the failure and replacing any faulty hardware.

If set to on, Postgres Pro will instead report an error but continue to run so that the data flushing operation can be retried in a later checkpoint. Only set it to on after investigating the operating system's treatment of buffered data in case of write-back failure.

`recovery_init_sync_method` (enum)

When set to `fsync`, which is the default, Postgres Pro will recursively open and synchronize all files in the data directory before crash recovery begins. The search for files will follow symbolic links for the WAL directory and each configured tablespace (but not any other symbolic links). This is intended to make sure that all WAL and data files are durably stored on disk before replaying changes. This applies whenever starting a database cluster that did not shut down cleanly, including copies created with `pg_basebackup`.

On Linux, `syncfs` may be used instead, to ask the operating system to synchronize the file systems that contain the data directory, the WAL files and each tablespace (but not any other file systems that may be reachable through symbolic links). This may be a lot faster than the `fsync` setting, because it doesn't need to open each file one by one. On the other hand, it may be slower if a file system is shared by other applications that modify a lot of files, since those files will also be written to disk. Furthermore, on versions of Linux before 5.8, I/O errors encountered while writing data to disk may not be reported to Postgres Pro, and relevant error messages may appear only in kernel logs.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`crash_info` (boolean)

When set to `on`, which is the default, Postgres Pro will write diagnostic information about a backend crash into a file.

This parameter can only be set at server start.

`crash_info_dump` (text)

Specifies a comma-separated list of character strings that contain data sources to provide data for a crash dump. Possible values of the strings are as follows:

- `queries` — query texts
- `memory_context` — memory context
- `system` — information on the OS
- `module` — information on modules loaded to the `postgres` process
- `cpuinfo` — information on the processor
- `virtual_memory` — information on virtual memory regions

The default value is `system,module,queries,memory_context`.

This parameter can only be set at server start.

`crash_info_location` (string)

Specifies the directory where information about a backend crash is to be stored. The value of `stderr` sends information about the crash to `stderr`. If this parameter is set to the empty string `''`, which is the default, the `$PGDATA/crash_info` directory is used. If you wish to keep the files elsewhere, create the target directory in advance and grant appropriate privileges.

This parameter can only be set at server start.

19.18. Preset Options

The following “parameters” are read-only. As such, they have been excluded from the sample `postgresql.conf` file. These options report various aspects of Postgres Pro behavior that might be of interest to certain applications, particularly administrative front-ends. Most of them are determined when Postgres Pro is compiled or when it is installed.

`block_size (integer)`

Reports the size of a disk block. It is determined by the value of `BLCKSZ` when building the server. The default value is 8192 bytes. The meaning of some configuration variables (such as [shared_buffers](#)) is influenced by `block_size`. See [Section 19.4](#) for information.

`data_checksums (boolean)`

Reports whether data checksums are enabled for this cluster. See [data checksums](#) for more information.

`data_directory_mode (integer)`

On Unix systems this parameter reports the permissions the data directory (defined by [data_directory](#)) had at server startup. (On Microsoft Windows this parameter will always display 0700.) See [group access](#) for more information.

`debug_assertions (boolean)`

Reports whether Postgres Pro has been built with assertions enabled. That is the case if the macro `USE_ASSERT_CHECKING` is defined when Postgres Pro is built (accomplished e.g., by the `configure` option `--enable-cassert`). By default Postgres Pro is built without assertions.

`huge_pages_status (enum)`

Reports the state of huge pages in the current instance: `on`, `off`, or `unknown` (if displayed with `postgres -C`). This parameter is useful to determine whether allocation of huge pages was successful under `huge_pages=try`. See [huge_pages](#) for more information.

`integer_datetimes (boolean)`

Reports whether Postgres Pro was built with support for 64-bit-integer dates and times. As of Postgres Pro 10, this is always `on`.

`in_hot_standby (boolean)`

Reports whether the server is currently in hot standby mode. When this is `on`, all transactions are forced to be read-only. Within a session, this can change only if the server is promoted to be primary. See [Section 26.4](#) for more information.

`lc_collate (string)`

Reports the locale in which sorting of textual data is done. See [Section 23.1](#) for more information. This value is determined when a database is created.

`lc_ctype (string)`

Reports the locale that determines character classifications. See [Section 23.1](#) for more information. This value is determined when a database is created. Ordinarily this will be the same as `lc_collate`, but for special applications it might be set differently.

`max_function_args (integer)`

Reports the maximum number of function arguments. It is determined by the value of `FUNC_MAX_ARGS` when building the server. The default value is 100 arguments.

`max_identifier_length (integer)`

Reports the maximum identifier length. It is determined as one less than the value of `NAMEDATALEN` when building the server. The default value of `NAMEDATALEN` is 64; therefore the default `max_identifier_length` is 63 bytes, which can be less than 63 characters when using multibyte encodings.

`max_index_keys (integer)`

Reports the maximum number of index keys. It is determined by the value of `INDEX_MAX_KEYS` when building the server. The default value is 32 keys.

`pgpro_build (string)`

Reports the commit ID of Postgres Pro source files.

`pgpro_edition (string)`

Reports the Postgres Pro edition as a string, i.e. `standard` or `enterprise`.

`pgpro_version (string)`

Reports the Postgres Pro server version as a string.

`segment_size (integer)`

Reports the number of blocks (pages) that can be stored within a file segment. It is determined by the value of `RELSEG_SIZE` when building the server. The maximum size of a segment file in bytes is equal to `segment_size` multiplied by `block_size`; by default this is 1GB.

`server_encoding (string)`

Reports the database encoding (character set). It is determined when the database is created. Ordinarily, clients need only be concerned with the value of [client_encoding](#).

`server_version (string)`

Reports the version number of the server. It is determined by the value of `PG_VERSION` when building the server.

`server_version_num (integer)`

Reports the version number of the server as an integer. It is determined by the value of `PG_VERSION_NUM` when building the server.

`shared_memory_size (integer)`

Reports the size of the main shared memory area, rounded up to the nearest megabyte.

`shared_memory_size_in_huge_pages (integer)`

Reports the number of huge pages that are needed for the main shared memory area based on the specified [huge_page_size](#). If huge pages are not supported, this will be `-1`.

This setting is supported only on Linux. It is always set to `-1` on other platforms. For more details about using huge pages on Linux, see [Section 18.4.5](#).

`ssl_library (string)`

Reports the name of the SSL library that this Postgres Pro server was built with (even if SSL is not currently configured or in use on this instance), for example `OpenSSL`, or an empty string if none.

`wal_block_size (integer)`

Reports the size of a WAL disk block. It is determined by the value of `XLOG_BLCKSZ` when building the server. The default value is 8192 bytes.

`wal_segment_size (integer)`

Reports the size of write ahead log segments. The default value is 16MB. See [Section 28.5](#) for more information.

19.19. Customized Options

This feature was designed to allow parameters not normally known to Postgres Pro to be added by add-on modules (such as procedural languages). This allows extension modules to be configured in the standard ways.

Custom options have two-part names: an extension name, then a dot, then the parameter name proper, much like qualified names in SQL. An example is `plpgsql.variable_conflict`.

Because custom options may need to be set in processes that have not loaded the relevant extension module, Postgres Pro will accept a setting for any two-part parameter name. Such variables are treated as placeholders and have no function until the module that defines them is loaded. When an extension module is loaded, it will add its variable definitions and convert any placeholder values according to those definitions. If there are any unrecognized placeholders that begin with its extension name, warnings are issued and those placeholders are removed.

19.20. Developer Options

The following parameters are intended for developer testing, and should never be used on a production database. However, some of them can be used to assist with the recovery of severely damaged databases. As such, they have been excluded from the sample `postgresql.conf` file. Note that many of these parameters require special source compilation flags to work at all.

`allow_in_place_tablespaces` (boolean)

Allows tablespaces to be created as directories inside `pg_tblspc`, when an empty location string is provided to the `CREATE TABLESPACE` command. This is intended to allow testing replication scenarios where primary and standby servers are running on the same machine. Such directories are likely to confuse backup tools that expect to find only symbolic links in that location. Only superusers and users with the appropriate `SET` privilege can change this setting.

`allow_system_table_mods` (boolean)

Allows modification of the structure of system tables as well as certain other risky actions on system tables. This is otherwise not allowed even for superusers. Ill-advised use of this setting can cause irretrievable data loss or seriously corrupt the database system. Only superusers and users with the appropriate `SET` privilege can change this setting.

`backtrace_functions` (string)

This parameter contains a comma-separated list of C function names. If an error is raised and the name of the internal C function where the error happens matches a value in the list, then a backtrace is written to the server log together with the error message. This can be used to debug specific areas of the source code.

Backtrace support is not available on all platforms, and the quality of the backtraces depends on compilation options.

Only superusers and users with the appropriate `SET` privilege can change this setting.

`debug_discard_caches` (integer)

When set to 1, each system catalog cache entry is invalidated at the first possible opportunity, whether or not anything that would render it invalid really occurred. Caching of system catalogs is effectively disabled as a result, so the server will run extremely slowly. Higher values run the cache invalidation recursively, which is even slower and only useful for testing the caching logic itself. The default value of 0 selects normal catalog caching behavior.

This parameter can be very helpful when trying to trigger hard-to-reproduce bugs involving concurrent catalog changes, but it is otherwise rarely needed.

This parameter is supported when `DISCARD_CACHES_ENABLED` was defined at compile time (which happens automatically when using the configure option `--enable-cassert`). In production builds, its value will always be 0 and attempts to set it to another value will raise an error.

`debug_io_direct` (string)

Ask the kernel to minimize caching effects for relation data and WAL files using `O_DIRECT` (most Unix-like systems), `F_NOCACHE` (macOS) or `FILE_FLAG_NO_BUFFERING` (Windows).

May be set to an empty string (the default) to disable use of direct I/O, or a comma-separated list of operations that should use direct I/O. The valid options are `data` for main data files, `wal` for WAL files, and `wal_init` for WAL files when being initially allocated.

Some operating systems and file systems do not support direct I/O, so non-default settings may be rejected at startup or cause errors.

Currently this feature reduces performance, and is intended for developer testing only.

`debug_parallel_query` (enum)

Allows the use of parallel queries for testing purposes even in cases where no performance benefit is expected. The allowed values of `debug_parallel_query` are `off` (use parallel mode only when it is expected to improve performance), `on` (force parallel query for all queries for which it is thought to be safe), and `regress` (like `on`, but with additional behavior changes as explained below).

More specifically, setting this value to `on` will add a `Gather` node to the top of any query plan for which this appears to be safe, so that the query runs inside of a parallel worker. Even when a parallel worker is not available or cannot be used, operations such as starting a subtransaction that would be prohibited in a parallel query context will be prohibited unless the planner believes that this will cause the query to fail. If failures or unexpected results occur when this option is set, some functions used by the query may need to be marked `PARALLEL UNSAFE` (or, possibly, `PARALLEL RESTRICTED`).

Setting this value to `regress` has all of the same effects as setting it to `on` plus some additional effects that are intended to facilitate automated regression testing. Normally, messages from a parallel worker include a context line indicating that, but a setting of `regress` suppresses this line so that the output is the same as in non-parallel execution. Also, the `Gather` nodes added to plans by this setting are hidden in `EXPLAIN` output so that the output matches what would be obtained if this setting were turned `off`.

`ignore_system_indexes` (boolean)

Ignore system indexes when reading system tables (but still update the indexes when modifying the tables). This is useful when recovering from damaged system indexes. This parameter cannot be changed after session start.

`post_auth_delay` (integer)

The amount of time to delay when a new server process is started, after it conducts the authentication procedure. This is intended to give developers an opportunity to attach to the server process with a debugger. If this value is specified without units, it is taken as seconds. A value of zero (the default) disables the delay. This parameter cannot be changed after session start.

`pre_auth_delay` (integer)

The amount of time to delay just after a new server process is forked, before it conducts the authentication procedure. This is intended to give developers an opportunity to attach to the server process with a debugger to trace down misbehavior in authentication. If this value is specified without units, it is taken as seconds. A value of zero (the default) disables the delay. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`trace_notify` (boolean)

Generates a great amount of debugging output for the `LISTEN` and `NOTIFY` commands. [client_min_messages](#) or [log_min_messages](#) must be `DEBUG1` or lower to send this output to the client or server logs, respectively.

`trace_sort` (boolean)

If `on`, emit information about resource usage during sort operations. This parameter is only available if the `TRACE_SORT` macro was defined when Postgres Pro was compiled. (However, `TRACE_SORT` is currently defined by default.)

`trace_locks (boolean)`

If on, emit information about lock usage. Information dumped includes the type of lock operation, the type of lock and the unique identifier of the object being locked or unlocked. Also included are bit masks for the lock types already granted on this object as well as for the lock types awaited on this object. For each lock type a count of the number of granted locks and waiting locks is also dumped as well as the totals. An example of the log file output is shown here:

```
LOG:  LockAcquire: new: lock(0xb7acd844) id(24688,24696,0,0,0,1)
      grantMask(0) req(0,0,0,0,0,0,0,0)=0 grant(0,0,0,0,0,0,0,0)=0
      wait(0) type(AccessShareLock)
LOG:  GrantLock: lock(0xb7acd844) id(24688,24696,0,0,0,1)
      grantMask(2) req(1,0,0,0,0,0,0,0)=1 grant(1,0,0,0,0,0,0,0)=1
      wait(0) type(AccessShareLock)
LOG:  UnGrantLock: updated: lock(0xb7acd844) id(24688,24696,0,0,0,1)
      grantMask(0) req(0,0,0,0,0,0,0,0)=0 grant(0,0,0,0,0,0,0,0)=0
      wait(0) type(AccessShareLock)
LOG:  CleanUpLock: deleting: lock(0xb7acd844) id(24688,24696,0,0,0,1)
      grantMask(0) req(0,0,0,0,0,0,0,0)=0 grant(0,0,0,0,0,0,0,0)=0
      wait(0) type(INVALID)
```

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`trace_lwlocks (boolean)`

If on, emit information about lightweight lock usage. Lightweight locks are intended primarily to provide mutual exclusion of access to shared-memory data structures.

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`trace_userlocks (boolean)`

If on, emit information about user lock usage. Output is the same as for `trace_locks`, only for advisory locks.

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`trace_lock_oidmin (integer)`

If set, do not trace locks for tables below this OID (used to avoid output on system tables).

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`trace_lock_table (integer)`

Unconditionally trace locks on this table (OID).

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`debug_deadlocks (boolean)`

If set, dumps information about all current locks when a deadlock timeout occurs.

This parameter is only available if the `LOCK_DEBUG` macro was defined when Postgres Pro was compiled.

`log_btree_build_stats (boolean)`

If set, logs system resource usage statistics (memory and CPU) on various B-tree operations.

This parameter is only available if the `BTREE_BUILD_STATS` macro was defined when Postgres Pro was compiled.

`wal_consistency_checking` (string)

This parameter is intended to be used to check for bugs in the WAL redo routines. When enabled, full-page images of any buffers modified in conjunction with the WAL record are added to the record. If the record is subsequently replayed, the system will first apply each record and then test whether the buffers modified by the record match the stored images. In certain cases (such as hint bits), minor variations are acceptable, and will be ignored. Any unexpected differences will result in a fatal error, terminating recovery.

The default value of this setting is the empty string, which disables the feature. It can be set to `all` to check all records, or to a comma-separated list of resource managers to check only records originating from those resource managers. Currently, the supported resource managers are `heap`, `heap2`, `btree`, `hash`, `gin`, `gist`, `sequence`, `spgist`, `brin`, and `generic`. Extensions may define additional resource managers. Only superusers and users with the appropriate `SET` privilege can change this setting.

`wal_debug` (boolean)

If on, emit WAL-related debugging output. This parameter is only available if the `WAL_DEBUG` macro was defined when Postgres Pro was compiled.

`ignore_checksum_failure` (boolean)

Only has effect if [data checksums](#) are enabled.

Detection of a checksum failure during a read normally causes Postgres Pro to report an error, aborting the current transaction. Setting `ignore_checksum_failure` to on causes the system to ignore the failure (but still report a warning), and continue processing. This behavior may *cause crashes, propagate or hide corruption, or other serious problems*. However, it may allow you to get past the error and retrieve undamaged tuples that might still be present in the table if the block header is still sane. If the header is corrupt an error will be reported even if this option is enabled. The default setting is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`zero_damaged_pages` (boolean)

Detection of a damaged page header normally causes Postgres Pro to report an error, aborting the current transaction. Setting `zero_damaged_pages` to on causes the system to instead report a warning, zero out the damaged page in memory, and continue processing. This behavior *will destroy data*, namely all the rows on the damaged page. However, it does allow you to get past the error and retrieve rows from any undamaged pages that might be present in the table. It is useful for recovering data if corruption has occurred due to a hardware or software error. You should generally not set this on until you have given up hope of recovering data from the damaged pages of a table. Zeroed-out pages are not forced to disk so it is recommended to recreate the table or the index before turning this parameter off again. The default setting is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`ignore_invalid_pages` (boolean)

If set to `off` (the default), detection of WAL records having references to invalid pages during recovery causes Postgres Pro to raise a PANIC-level error, aborting the recovery. Setting `ignore_invalid_pages` to on causes the system to ignore invalid page references in WAL records (but still report a warning), and continue the recovery. This behavior may *cause crashes, data loss, propagate or hide corruption, or other serious problems*. However, it may allow you to get past the PANIC-level error, to finish the recovery, and to cause the server to start up. The parameter can only be set at server start. It only has effect during recovery or in standby mode.

`jit_debugging_support` (boolean)

If LLVM has the required functionality, register generated functions with GDB. This makes debugging easier. The default setting is `off`. This parameter can only be set at server start.

`jit_dump_bitcode` (boolean)

Writes the generated LLVM IR out to the file system, inside [data_directory](#). This is only useful for working on the internals of the JIT implementation. The default setting is `off`. Only superusers and users with the appropriate `SET` privilege can change this setting.

`jit_expressions` (boolean)

Determines whether expressions are JIT compiled, when JIT compilation is activated (see [Section 30.2](#)). The default is `on`.

`jit_profiling_support` (boolean)

If LLVM has the required functionality, emit the data needed to allow `perf` to profile functions generated by JIT. This writes out files to `~/.debug/jit/`; the user is responsible for performing cleanup when desired. The default setting is `off`. This parameter can only be set at server start.

`jit_tuple_deforming` (boolean)

Determines whether tuple deforming is JIT compiled, when JIT compilation is activated (see [Section 30.2](#)). The default is `on`.

`remove_temp_files_after_crash` (boolean)

When set to `on`, which is the default, Postgres Pro will automatically remove temporary files after a backend crash. If disabled, the files will be retained and may be used for debugging, for example. Repeated crashes may however result in accumulation of useless files. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`send_abort_for_crash` (boolean)

By default, after a backend crash the postmaster will stop remaining child processes by sending them `SIGQUIT` signals, which permits them to exit more-or-less gracefully. When this option is set to `on`, `SIGABRT` is sent instead. That normally results in production of a core dump file for each such child process. This can be handy for investigating the states of other processes after a crash. It can also consume lots of disk space in the event of repeated crashes, so do not enable this on systems you are not monitoring carefully. Beware that no support exists for cleaning up the core file(s) automatically. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`send_abort_for_kill` (boolean)

By default, after attempting to stop a child process with `SIGQUIT`, the postmaster will wait five seconds and then send `SIGKILL` to force immediate termination. When this option is set to `on`, `SIGABRT` is sent instead of `SIGKILL`. That normally results in production of a core dump file for each such child process. This can be handy for investigating the states of “stuck” child processes. It can also consume lots of disk space in the event of repeated crashes, so do not enable this on systems you are not monitoring carefully. Beware that no support exists for cleaning up the core file(s) automatically. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`debug_logical_replication_streaming` (enum)

The allowed values are `buffered` and `immediate`. The default is `buffered`. This parameter is intended to be used to test logical decoding and replication of large transactions. The effect of `debug_logical_replication_streaming` is different for the publisher and subscriber:

On the publisher side, `debug_logical_replication_streaming` allows streaming or serializing changes immediately in logical decoding. When set to `immediate`, stream each change if the [streaming](#) option of `CREATE SUBSCRIPTION` is enabled, otherwise, serialize each change. When set to `buffered`, the decoding will stream or serialize changes when `logical_decoding_work_mem` is reached.

On the subscriber side, if the `streaming` option is set to `parallel`, `debug_logical_replication_streaming` can be used to direct the leader apply worker to send changes to the shared mem-

ory queue or to serialize all changes to the file. When set to `buffered`, the leader sends changes to parallel apply workers via a shared memory queue. When set to `immediate`, the leader serializes all changes to files and notifies the parallel apply workers to read and apply them at the end of the transaction.

19.21. Distributed Cluster Configuration

`shardmanctl` uses the `sdmspec.json` configuration file during Postgres Pro Shardman cluster initialization. A `shardman-utils` package provides a sample configuration file.

`sdmspec.json` file contains basic filesystem paths used by Postgres Pro Shardman, global settings of the cluster, database-related settings, i. e., administrative and replication user logins and authentication method, FDW parameters and shard configuration (`ShardSpec`).

Note that there is a number of the internal Postgres Pro Shardman parameters that, if modified by user, can result in the total cluster failure. These parameters are:

- `shardman.cluster_uuid` defines the version of a running cluster that the node belongs to.
- `shardman.config_uuid` defines the config version. Ignored if set via `shardmanctl config update` or `shardmanctl init`.
- `shardman.manual_execution` controls the consistent work with the global objects.
- `shardman.silk_never_restart` prohibits the multiplexer workers restart in case of an error.
- `shardman.pre_promote_mode` applies the consistent promotion mechanism (from standby to primary).

19.21.1. List of Parameters

Repfactor

Integer determining how many replicas `shardmanctl` should configure for each DBMS. This setting can only be changed for a Postgres Pro Shardman cluster with a manual-topology mode.

DataDir

Allows you to specify a directory other than the default one (`/var/lib/pgpro/sdm-17/data`) for storing data. This parameter cannot be changed after the cluster has been initialized.

PGsInitialPort

Ports starting with this integer are assigned to PostgreSQL instances. This parameter cannot be changed after the cluster has been initialized.

SilkInitialPort

Ports starting with this integer are assigned to Silk (Postgres Pro Shardman InterLink) instances. This parameter cannot be changed after the cluster has been initialized.

AuthMethod

Authentication method used by the administrative user to connect to the DBMS. Can be any authentication method supported by PostgreSQL. `scram-sha-256` is currently recommended. `md5` is currently allowed but not recommended. This parameter cannot be changed after the cluster has been initialized. Located under a separate `Users` block for each array element.

Default: `trust`.

Groups

An array that can have two possible values, `su` for superuser or `repl` for replication.

HTTP

Defines settings for the secure HTTP/HTTPS connection, with *Port* being an API port, and *PortMetrics* being a port for the metrics. If these ports are the same, then API and metrics listen to the same port.

Default: 15432.

Name

Name of the user. Created on cluster initialization. Defaults to the name of the effective user running `shardmanctl init`. This parameter cannot be changed after the cluster has been initialized. Located under a separate `Users` block for each array element.

Password

Password for the user. Can be changed using `shardmanctl config update credentials`. Located under a separate `Users` block for each array element.

PgSuSSLCert

Client certificate for the administrative DBMS user.

PgSSLRootCert

Location of the root certificate file for the DBMS user connection.

PgSuSSLKey

Client private key for the administrative DBMS user.

PgSSLMode

SSL mode for the DBMS user. Allowed values: `verify-ca` and `verify-full`.

PgReplSSLCert

Client certificate for the replication DBMS user.

PgReplSSLKey

Client private key for the replication DBMS user.

ShardSpec

Shard cluster specification. For more details, see [ShardSpec Parameters](#). Can be changed using `shardmanctl config update`.

FDWOptions

This object contains FDW settings.

These settings can be changed using `shardmanctl config update` (with the exception of settings related to authorization, server connection, SSL and Kerberos, as well as the `service`, `target_session_attrs` options).

Foreign servers corresponding to Postgres Pro Shardman replication groups will also get `extended_features` setting automatically enabled. Never set this parameter for postgres fdw foreign servers which you define for your own purposes (for example, to load data into Postgres Pro Shardman cluster).

19.21.2. ShardSpec Parameters

The `ShardSpec` specification lists a set of parameters that can be tuned for a specific replication group.

pgHBA

JSON array of `pg_hba.conf` strings. The default value allows user from the `su` group access from anywhere with `AuthMethod` authentication method. If the value of `defaultSUReplAccessMode` is `strict`, `pg_hba.conf` strings must explicitly allow users from the groups `su` or `repl` access from all Postgres Pro Shardman cluster nodes.

Cluster software does not fail even if writing to `pg_hba.conf` is forbidden.

forceSuUserLocalPeerAuth

When enabled, it sets a peer authentication via unix socket for the `postgres` user, if `strictUserHBA` is not set to `true`.

Default: `false`.

maxSynchronousStandbys

Maximum number of required synchronous standbys when synchronous replication is enabled. Should be \geq [Repfactor](#) in a Postgres Pro Shardman cluster. Default: `Repfactor`.

strictUserHBA

Prohibits adding automatically generated lines to `pg_hba.conf` file. Default: `false`.

automaticPgRestart

Determines whether a DBMS instance should be automatically restarted after a change of the [pgParameters](#) hash table that requires a restart. Should be enabled in a Postgres Pro Shardman cluster.

Default: `true`.

masterDemotionEnabled

Enable master demotion in case the replica group master has lost connectivity with `etcd`. The master attempts to connect to each of its standby nodes to determine if any of them has become the master. If it discovers another master, it shuts down its own DBMS instance until the connectivity with `etcd` is restored. If the master fails to connect to one of its standby nodes for a long time, a DBMS instance shutdown occurs.

Default: `false`.

masterDemotionTimeout

The timeout during which the master attempts to connect to its standbys in cases where connectivity with `etcd` is lost. Works only if the `masterDemotionEnabled` parameter is set to `true`.

Default: `30s`.

minSyncMonitorEnabled

Enable the monitor for the `MinSynchronousStandbys` value for every replica group. If a node loses connection with the cluster (all keepers are unhealthy: a keeper does not update its state longer than `minSyncMonitorUnhealthyTimeout`), the monitor decreases the `MinSynchronousStandbys` value for every replica group related to the disconnected node to the maximum available value. This allows preventing the read-only condition caused by the fake replica. The maximum available value is always less than or equal to the value specified in the cluster configuration. If all keepers related to the disconnected node become healthy, the monitor changes `MinSynchronousStandbys` value for the replica group to the value specified in the cluster configuration.

Default: `false`.

minSyncMonitorUnhealthyTimeout

Time interval after which the node (and all keepers related to this node) will be considered in an unhealthy condition. Works only if the `minSyncMonitorEnabled` parameter is set to `true`.

Default: 30s.

`syncPointMonitorEnabled`

Enable the monitor that creates a syncpoint every minute, ensuring the Postgres Pro Shardman can restore to a consistent LSN. At each syncpoint, the cluster's state is consistent, meaning that all transactions are complete. If this parameter is set to `true`, PITR will be guaranteed to work. `syncPointMonitorEnabled` saves the syncpoint history in etcd with the `shardman/{cluster_name}/data/cluster/syncpoints` key.

Default: `false`.

`syncPointMonitorLimit`

Specifies the limit for the size of the syncpoint history saved by `syncPointMonitorEnabled`.

Default: 60.

`shardmanSchemaDumpMonitorDisabled`

Disables the `shardmanSchemaDumpMonitorDisabled` monitor that periodically creates a dump of the shardman schema.

Default: `true`.

`shardmanSchemaDumpMonitorTick`

Sets how often the `shardmanSchemaDumpMonitor` monitor creates a dump.

Default: every 20 minutes.

`pgParameters`

Hash table that determines PostgreSQL settings, including [settings for a distributed system](#).

`biHAReplicaWithProbackup`

If enabled, new BiHA followers and `referee-with-wal` will use `probackup` as synchronization method.

19.21.3. PostgreSQL Settings for a Distributed System

The following settings in `pgParameters` are specific for a distributed system:

`enable_sql_func_custom_plans` (boolean)

If enabled, custom plans can be created to execute statements inside SQL functions. These plans depend on the parameter values.

Query plans can be cached within one query. First, the plan is built five times with different parameter values, then a generic plan is created regardless of the values. If custom and generic plan price is slightly different, then the generic plan is cached and is set to be used in the future. However, custom plans allow a more effective way of excluding queries to the sharded table partitions if the choice of these partitions depends on the query option.

Default: `off`.

`enable_merge_append` (boolean)

Enables the use of `MergeAppend` plans by the query planner.

Default: `on`.

`enable_async_merge_append` (boolean)

Enables or disables the query planner's use of `async-aware merge append` plan types. The default is `on`.

`csn_snapshot_defer_time` (integer)

Specifies the minimal age of records that are allowed to be vacuumed, in seconds.

For distributed system specifics, see [csn_snapshot_defer_time](#).

All global transactions must start on all participant nodes within `csn_snapshot_defer_time` seconds after start, otherwise, they are aborted with a “`csn snapshot too old`” error.

Default: 15.

`csn_commit_delay` (integer)

Specifies the maximum possible clock skew (in nanoseconds) in the cluster. Adds a delay before every commit in the system to ensure external consistency. If set to 0, external consistency is not guaranteed. Value suffixes `ns`, `us`, `ms` and `s` are allowed.

Default: 0.

`csn_lsn_map_size` (integer)

Size of CSNLSNMap.

The commit record of each completed transaction in Postgres Pro Shardman contains the assigned CSN for this transaction. This value, together with the LSN of this record, forms a pair of values (CSN, LSN). Each of the cluster nodes stores a certain number of such pairs in RAM in a special structure - the CSNLSNMap. This map is used to get the syncpoint. See the “Syncpoints and Consistent Backup” section of the [Internals](#) chapter for more information.

Default: 1024.

`csn_max_shift_error` (boolean)

When checked against the `csn_max_shift` value, raises an error if the `csn_max_shift` value is exceeded.

Default: off.

`csn_max_shift` (integer)

Maximum CSN shift in seconds for distributed queries and imported snapshots. If the shift exceeds the `csn_max_shift` value, an error or warning will occur. If the value is set to 0, no check is run.

Default: 15 seconds.

`foreign_analyze_interval` (integer)

Specifies how often foreign statistics should be gathered during autovacuum, in seconds. If the value of `foreign_analyze_interval` is less than [autovacuum_naptime](#), foreign statistics will be gathered each `autovacuum_naptime` seconds.

Default: 60.

`foreign_join_fast_path` (boolean)

Turns on a fast path for foreign join planning. When it is on, foreign join paths for `SELECT` queries are searched before all other possible paths and the search stops for a join as soon as a foreign join path is found.

Default: off.

`optimize_correlated_subqueries` (boolean)

Enables or disables the query planner's logic of transforming correlated subqueries into semi-joins.

Default: on.

`port` (integer)

A TCP port the server listens on. For a Postgres Pro Shardman cluster, the `port` is assigned automatically by the system and is based on the `PGsInitialPort` parameter. If changed manually, the value will be overwritten by the configuration parameter that is automatically assigned.

`enable_partition_pruning_extra` (boolean)

Enables the extended partition pruning for the prepared queries with a known partitioning key. If turned on, the partition-wise join plans can be pruned.

Default: `off`.

`postgres_fdw.enforce_foreign_join` (boolean)

Turns on alternative estimations for foreign join costs, which highly increases chances for join of several foreign tables referring to the same server to be pushed down. The cost of original join is estimated as $(1 - 1/(cost + 1))$, where `cost` is an originally estimated cost for this remote join.

Default: `off`.

`postgres_fdw.foreign_explain` (enum)

Defines how to include the `EXPLAIN` command output from the remote servers if the query plan contains `ForeignScan` nodes. The possible values are: `none` to exclude the `EXPLAIN` output from the remote servers, `full` to include the `EXPLAIN` output from the remote servers, `collapsed` to include the `EXPLAIN` output only for the first `ForeignScan` node under its `Append/MergeAppend`.

Default: `collapsed`.

`postgres_fdw.use_twophase` (enum)

Sets `postgres_fdw` to use the two-phase commit (2PC) protocol for distributed transactions.

This `postgres_fdw` parameter forces it to use a two-phase commit if the transaction touches several nodes. When set to `auto`, a two-phase commit is only used in transactions with `enable_csn_snapshot=true` and isolation level equal to or higher than `REPEATABLE READ`.

Temporary tables cannot be used in 2PC transactions.

Default: `auto`.

`postgres_fdw.additional_ordered_paths` (boolean)

When enabled, sorting on the remote server is considered if it allows performing `MergeJoin` or `MergeAppend` operations. This parameter is enabled by default in new installations but must be explicitly enabled in upgraded clusters.

`shardman.broadcast_ddl` (boolean)

Sets Postgres Pro Shardman extension to broadcast DDL statements to all replication groups.

When this parameter is on, Postgres Pro Shardman extension broadcasts supported DDL statements to all replication groups if it does make sense for those statements. You can enable/disable this behavior anytime. This parameter is not honored when set in configuration file.

Default: `off`.

`shardman.enable_limit_pushdown` (boolean)

Enable pushing down limit clauses through the underlying appends. When on, Postgres Pro Shardman optimizer will try to push down a limit clause to the subpaths of the underlying `Append/MergeAppend` plan node if they reference `postgres_fdw` foreign tables. This optimization works only for `SELECT` plans when limit option is represented as a constant or a parameter. It is also restricted for `Append`

paths, corresponding to a partitioned table. The optimization does not work for `SELECT` with locking clauses (`SELECT FOR UPDATE/NO KEY UPDATE/FOR SHARE/KEY SHARE`).

Default: `on`.

`shardman.num_parts` (integer)

Specifies the default number of sharded table partitions.

A sharded table has this default number of partitions unless `num_parts` is specified in [CREATE TABLE](#).

To allow scaling, `shardman.num_parts` should be larger than the expected maximum number of nodes in a Postgres Pro Shardman cluster.

Possible values are from 1 to 1000.

Default: 20.

`shardman.rgid` (integer)

Specifies the replication group ID of a Postgres Pro Shardman node.

This parameter is set by Postgres Pro Shardman utilities when the node is added to the cluster and should never be changed manually.

Default: -1.

`shardman.sync_schema` (boolean)

Sets Postgres Pro Shardman to propagate all DDL statements that touch sharded and global relations to all replication groups.

When this parameter is on, Postgres Pro Shardman broadcasts all supported utility statements touching sharded and global relations to all replication groups. It is not recommended to turn this off. This parameter is not honored when set in configuration file.

Default: `on`.

`shardman.sync_cluster_settings` (boolean)

Enables cluster-wide synchronization of configuration parameters set by user. The configuration parameters are propagated with each remote query.

Default: `on`.

`shardman.sync_cluster_settings_blacklist` (boolean)

Excludes the options not to be propagated to a remote cluster.

Default: local system configuration parameters that are never synchronized.

`shardman.query_engine_mode` (enum)

Switches between modes of query planning/execution. Possible values are `none` and `text`.

`none` means that query planning/execution will not use the Silk transport.

`text` means that the text query representation is transferred via Silk transport for remote execution.

Default: `none`.

`shardman.silk_use_ip` (string)

Silk transport uses IP address specified by this parameter for node identification. If the host name is specified, it is resolved and the first IP address corresponding to this name, is used.

Default: node hostname.

`shardman.silk_listen_ip (string)`

The Silk routing daemon listens for incoming connections on this IP address. If the host name is specified, it is resolved and the first IP address corresponding to this name, is used.

Default: node hostname.

`shardman.silk_use_port (integer)`

The Silk routing daemon listens for incoming connections on this port. This setting should be the same for all nodes in the Postgres Pro Shardman cluster.

Default: 8888.

`shardman.silk_tracepoints (bool)`

Enables tracing of queries passing through the Silk pipeline. The tracing results can be accessed by running the `EXPLAIN` command with `ANALYZE` set to `ON`.

Default: off.

`shardman.silk_num_workers (integer)`

Number of background workers allocated for distributed execution. This setting must be less than `max_worker_processes` (including auxiliary postgres worker processes).

Default: 2.

`shardman.silk_stream_work_mem (integer)`

Sets the base maximum amount of memory to be used by a Silk stream (as a buffer size) before writing to the temporary disk files. If this value is specified without units, the default is kilobytes.

Note that most queries can perform multiple fetch operations at the same time, usually one for each remote partition of a sharded table, if any. Each fetch operation is generally allowed to use as much memory as this value specifies before it starts to write data into temporary files. Also, several running sessions can execute such operations concurrently. Therefore, the total memory used by Silk for buffers could be many times the value of `shardman.silk_stream_work_mem` and is correlated with [shardman.num_parts](#). Thus, mind this fact when choosing the value.

Default: 16MB.

`shardman.silkworm_fetch_size (integer)`

Number of rows in a chunk that the `silkworm` worker extracts and sends to the multiplexer as a result, per one reading iteration.

Default: 100.

`shardman.silk_unassigned_job_queue_size (integer)`

Size of queue for jobs that have not yet been assigned to the `silkworm` multiplexer workers, in case all the workers are busy.

Default: 1024.

`shardman.silk_max_message (integer)`

Maximum message size that can be transferred with Silk, in bytes. Note that this parameter does not limit the maximum size of the result returned by the query. It only affects messages sent to workers. Increasing this parameter value will result in a proportional memory increase consumed by Postgres Pro Shardman. It is strongly recommended to use the default value unless there is an urgent need.

Default: 524288.

`shardman.silk_hello_timeout` (integer)

Handshake timeout between multiplexers of different nodes, in seconds.

Default: 3.

`shardman.silk_scheduler_mode` (enum)

Enables additional CPU scheduling settings for multiplexer processes (`silkroad` and `silkworm`).

When this parameter is `fifo`, Postgres Pro Shardman assigns scheduling policy `SCHED_FIFO` for processes `silkroad` and each of `silkworm`. It assigns the static scheduling priority (`sched_priority`) to values `shardman.silkroad_sched_priority` and `shardman.silkworm_sched_priority` respectively.

This setting improves silk transport performance while it operates under heavy CPU load.

Note that `postgres` binary need to have `CAP_SYS_NICE` capability to use this option. If no appropriate capability was assigned to the process, enabling this setting will have no effect. The capability must be assigned to `postgres` binary before starting `postgres`. `postgres` (i.e. processes `silkroad` and `silkworm`) will apply scheduling options once during service start. You need restart `postgres` service if you want to change scheduling options.

Default: `none`.

To set capability you need execute following command once after `postgres` installed:

```
$ sudo setcap cap_sys_nice+ep /opt/pgpro/sdm-17/bin/
postgres
```

Replace `/opt/pgpro/sdm-17/bin/postgres` to the correct path to your `postgres` binary if needed. Also note that your filesystem should support extended file attributes. You need set this for each node in the cluster to take the full effect.

In the Linux kernel, there is a mechanism called real-time throttling, which is designed to prevent tasks with real-time scheduling policies (like `SCHED_FIFO`) from monopolizing CPU resources. This ensures that other tasks with lower priorities, typically scheduled under the `SCHED_OTHER` policy, still get some amount of the CPU time. This mechanism is controlled by two parameters, exported into the `proc` filesystem or the `sysctl` interface:

- `/proc/sys/kernel/sched_rt_period_us` sets the duration of a scheduling period in microseconds. During this period, both real-time and non-real-time tasks share CPU time.
- `/proc/sys/kernel/sched_rt_runtime_us` specifies how much of the scheduling period is allocated to real-time tasks (with `SCHED_FIFO`). The remainder of the time is left for non-real-time tasks (`SCHED_OTHER`).

A typical and acceptable configuration for Postgres Pro Shardman might set these parameters as follows:

```
# cat /proc/sys/kernel/sched_rt_period_us
1000000
# cat /proc/sys/kernel/sched_rt_runtime_us
950000
```

This configuration allows real-time tasks to use up to 950 milliseconds of each second, leaving 50 milliseconds for non-real-time tasks.

However, in some Linux distributions, the default values for these parameters might be set so low (or even to zero) that real-time tasks receive very little or no CPU time. This can make real-time scheduling ineffective or prevent the configuration from being applied. For example, attempting to manually set a task to the `SCHED_FIFO` priority using `chrt` might result in an error like:

```
$ sudo chrt -f -p 2 $(pgrep -f silkroad)
chrt: failed to set pid 1897706's policy: Operation not permitted
```

This error indicates that the kernel parameters are not configured correctly. In such cases, run the following:

```
echo 1000000 > /proc/sys/kernel/sched_rt_period_us
echo 950000 > /proc/sys/kernel/sched_rt_runtime_us
```

Or add the corresponding values into `/etc/sysctl.conf` and reload the settings using `sysctl -p`:

```
kernel.sched_rt_period_us = 1000000
kernel.sched_rt_runtime_us = 950000
```

`shardman.silkroad_sched_priority` (integer)

Value of static scheduling priority (`sched_priority`) for `silkroad` process. It only makes sense if `shardman.silk_scheduler_mode` equals to 'fifo'.

Default: 2.

`shardman.silkworm_sched_priority` (integer)

Value of static scheduling priority (`sched_priority`) for `silkworm` processes (the same value for each of them). It only makes sense if `shardman.silk_scheduler_mode` equals to 'fifo'.

Default: 1.

`shardman.silk_set_affinity` (bool)

Enables pinning of multiplexer processes (`silkroad` and `silkworm`) to CPU cores to eliminate negative effects of thread's cross-cpu migration.

When this parameter is `true`, `silkroad` process will be pinned to the first available CPU core and `silkworm` processes (all of them) will be pinned to all available CPU cores except the first one.

This setting improves silk transport performance while it operates under heavy CPU load.

Note that `postgres` binary need to have `CAP_SYS_NICE` capability to use this option. If no appropriate capability was assigned to the process, enabling this setting will have no effect. The capability must be assigned to `postgres` binary before starting `postgres` (i.e. processes `silkroad` and `silkworm`) will apply affinity options once during service start. You need restart `postgres` service if you want to change affinity options.

To set capability you need execute following command once after `postgres` installed:

```
$ sudo setcap cap_sys_nice+ep /opt/pgpro/sdm-17/bin/
postgres
```

Replace `/opt/pgpro/sdm-17/bin/postgres` to the correct path to your `postgres` binary if needed. Also note that your filesystem should support extended file attributes. You need set this for each node in the cluster to take the full effect.

Default: false.

`shardman.silk_flow_control` (boolean)

Controls the mode of handling read events. It has three possible values: `none`, `round_robin`, and `shortest_job_first`.

The `none` mode means no control nor additional overhead. Yet in this case, the channel may become occupied by just one distributed query.

The `round_robin` mode means the events created earlier are the first ones to be processed, for each event loop. If enabled, all the backends are grouped, and the client backends are prioritized over the other.

The `shortest_job_first` mode means full control over the traffic. If enabled, all the backends are grouped, and the client backends are prioritized over the others, along with the workers with the least session traffic.

Default: `round_robin`.

`shardman.silk_track_time` (boolean)

Enables or disables the metrics with prefix `transferred_` and time-based metrics (with prefixes `read_efd_`, `write_efd_`, and `sort_time_`). If disabled, these metrics have 0 values.

Default: `off`.

`shardman.silk_tracelog` (bool)

Enables or disables Silk logging.

Default: `off`.

`shardman.silk_tracelog_category` (string)

Defines the Silk message categories to be traced.

Default: `streams, routing, events`.

`shardman.database` (string)

Name of the database that all Silk workers connect to.

Default: `postgres`.

`shardman.monitor_interval` (integer)

`shardman.monitor_interval` is deprecated and acts as noop.

Use `shardman.monitor_dxact_interval` instead.

`shardman.monitor_dxact_interval` (integer)

Interval between checks for outdated prepared transactions.

The Postgres Pro Shardman monitor background process wakes up every `shardman.monitor_dxact_interval` seconds and attempts to check and resolve any prepared transactions that did not complete and became outdated for some reason. To resolve these transactions, the Postgres Pro Shardman monitor process determines the coordinator of the transaction and requests the transaction status from the coordinator. Based on the status of the transaction, Postgres Pro Shardman monitor will either roll back or commit the transaction.

To disable the prepared transaction resolution logic, set `shardman.monitor_dxact_interval` to 0.

Default: 5 (seconds).

`shardman.monitor_trim_csnxid_map_interval` (integer)

Each cluster node freezes its own `xmin` value for `csn_snapshot_defer_time` seconds to support global transactions. Large `csn_snapshot_defer_time` values can negatively impact the performance. Postgres Pro Shardman monitor has a routine that every `shardman.monitor_trim_csnxid_map_interval` seconds updates `xmin` on all nodes to the minimum possible value (taking into account active transactions).

The background routine will run on only one node in the Postgres Pro Shardman cluster. Note that this will give an additional load on this node.

To disable such updates, set `shardman.monitor_trim_csnxid_map_interval` to 0.

Default: 5 (seconds).

`shardman.monitor_dxact_timeout` (integer)

Maximum allowed age of prepared transactions before a resolution attempt.

During the resolution of a [prepared transaction](#), Postgres Pro Shardman monitor determines whether the transaction is outdated or not. A transaction becomes outdated if it was prepared more than `shardman.monitor_dxact_timeout` seconds ago.

Default: 5 (seconds).

`shardman.trim_csnxid_map_naptime` (integer)

Specifies the minimum delay between `xmin` updates on all nodes. See [shardman.monitor_trim_csnxid_map_interval](#) for more information.

Possible values are from 1 to 600.

Default: 5.

`shardman.monitor_deadlock_interval` (integer)

Interval between checks for distributed deadlock conditions.

The Postgres Pro Shardman monitor background process wakes up every `shardman.monitor_deadlock_interval` seconds and searches for distributed deadlocks in the cluster. It gathers information about mutual locks from all nodes and looks for circular dependencies between transactions. If it detects a deadlock, it resolves it by canceling one of the backend processes involved in the lock.

To disable the distributed deadlock resolution logic, set `shardman.monitor_deadlock_interval` to 0.

Default: 2 (seconds).

`postgres_fdw.remote_plan_cache` (boolean) — EXPERIMENTAL

Enables remote plan caching for FDW queries produced by locally cached plans.

Default: off.

`shardman.plan_cache_mem` (integer) — EXPERIMENTAL

Specifies how much memory per worker can be used for remote plan caches.

Default: 0 (caches are disabled).

`shardman.gt_batch_size` (integer) —

Specifies the buffer size for `INSERT` and `DELETE` commands executed on a global table.

Default: 64K.

`track_fdw_wait_timing` (boolean)

The statistics for the network latency (wait time) for inter-cluster operations, in milliseconds. It can be accessed by running the `EXPLAIN` command with the `network` parameter enabled.

Default: on.

`track_xact_time` (boolean)

Enables or disables statistics collection for time spent on a transaction.

Default: off.

`enable_non_equivalence_filters` (boolean)

Enables the optimizer to generate additional non-equivalence conditions using equivalence classes.

Default: off.

`optimize_row_in_expr` (boolean)

Enables the optimizer to generate additional conditions from the `IN ()` expression.

Default: off.

19.21.4. Examples

19.21.4.1. Spec File for a Cluster with Enabled `scram-sha-256` Authentication

Note

The initial configuration file should be generated with the following command:

```
shardmanctl config generate > sdmspec.json
```

The example below is for educational purposes only and may lack the latest updates.

This is the contents of an example `sdmspec.json` configuration file:

```
{
  "ConfigVersion": "1",
  "Repfactor": 1,
  "PGsInitialPort": 5432,
  "SilkInitialPort": 8000,
  "BiHAInitialPort" : 25432,
  "HTTP": {
    "Port": 15432,
    "PortMetrics": 15432
  },
  "Users": [
    {
      "Name": "postgres",
      "Groups": [ "su"],
      "AuthMethod": "scram-sha-256",
      "Password": "changeMe"
    },
  ],
}
```

```

    "Name": "biha_replication_user",
    "Groups": ["repl"],
    "AuthMethod": "scram-sha-256",
    "Password": "changeMe"
  },
],
"ShardSpec": {
  "usePgrewind": true,
  "pgParameters": {
    "csn_snapshot_defer_time": "300",
    "enable_csn_snapshot": "on",
    "enable_csn_wal": "true",
    "shardman.query_engine_mode": "text",
    "shardman.silk_num_workers": "8",
    "max_connections": "600",
    "max_files_per_process": "65535",
    "max_logical_replication_workers": "14",
    "max_prepared_transactions": "200",
    "max_worker_processes": "24",
    "shared_preload_libraries": "postgres_fdw, shardman, biha"
  },
  "pgHBA": [
    "host replication postgres 0.0.0.0/0 scram-sha-256",
    "host replication postgres ::0/0 scram-sha-256"
  ],
  "automaticPgRestart": true,
  "masterDemotionEnabled": false
},
"FDWOptions": {
  "async_capable": "on",
  "batch_size": "100",
  "connect_timeout": "5",
  "fdw_tuple_cost": "0.2",
  "fetch_size": "50000",
  "tcp_user_timeout": "10000"
}
}

```

From that configuration file, you can see that a Postgres Pro Shardman cluster initialized with this spec file has `Repfactor` equal to 1 (one replica for each master). The configuration file also shows that two special users are created in this cluster — superuser `postgres` and replication user `repluser` with `ChangeMe` passwords. They can be authenticated using the `md5` or `scram-sha-256` authorization method. One `postgres_fdw` fetch operation will get up to 50000 rows from the remote server. The cost of fetching one row is set to a reasonably high value to make PostgreSQL planner consider conditions pushdown-attractive. `pg_hba.conf` settings allow `postgres` user access from anywhere using a replication protocol; all other users can access any database from anywhere. Since `defaultSUReplAccessMode` is not set to `strict`, utilities will automatically add entries that allow `PgSuUsername` user's (`postgres`) access to any database from anywhere and `PgReplUsername` user's (`repluser`) replication access from anywhere.

Several important parameters for a distributed system are set in the `pgParameters` hash table. For regular system types see [Section 19.7](#). These are:

`wal_level`

Should be set to `logical` for Postgres Pro Shardman to work correctly.

`shared_preload_libraries`

Should include `postgres_fdw` and `shardman` extensions in the specified order.

`max_logical_replication_workers`

Should be rather high since the rebalance process uses up to `max(max_replication_slots, max_logical_replication_workers, max_worker_processes, max_wal_senders)/3` concurrent threads.

`max_prepared_transactions`

Should be rather high since Postgres Pro Shardman utilities use the 2PC protocol. If `postgres_fdw.use_twophase` is `true`, `postgres_fdw` also uses 2PC.

`enable_csn_snapshot`

Should be enabled to achieve a true `REPEATABLE READ` isolation level in a distributed system.

`csn_snapshot_defer_time`

All global transactions must start on all participant nodes within `csn_snapshot_defer_time` seconds after start, otherwise they will be aborted.

`enable_partitionwise_aggregate`

`enable_partitionwise_join`

Set to `on` to enable optimizations for partitioned tables.

19.21.4.2. Spec File for a Cluster with Enabled Certificate Authentication

This is the contents of an example `sdmspec.json` configuration file:

```
{
  "ConfigVersion": "1",
  "HTTP": {
    "Port": 15432,
    "PortMetrics": 15432
    "SSLKey": "/pgpro/ssl/server.key",
    "SSLCert": "/pgpro/ssl/server.crt"
  },
  "Users": [
    {
      "Name": "postgres",
      "SSLKey": "/var/lib/postgresql/.ssh/client.key",
      "SSLCert": "/var/lib/postgresql/.ssh/client.crt",
      "Groups": ["su"],
      "AuthMethod": "scram-sha-256"
    },
    {
      "Name": "biha_replication_user",
      "SSLKey": "/var/lib/postgresql/.ssh/repluser.key",
      "SSLCert": "/var/lib/postgresql/.ssh/repluser.crt",
      "Groups": ["repl"],
      "AuthMethod": "scram-sha-256"
    }
  ],
  "ShardSpec": {
    "synchronousReplication": true,
    "usePgrewind": true,
    "certVersion": true,
    "pgParameters": {
      "ssl": "on",
      "ssl_cert_file": "/var/lib/postgresql/.ssh/server.crt",
      "ssl_key_file": "/var/lib/postgresql/.ssh/server.key",

```

```

    "ssl_ca_file": "/var/lib/postgresql/.ssh/ca.crt",
    "csn_snapshot_defer_time": "300",
    "enable_csn_snapshot": "on",
    "enable_csn_wal": "true",
    "log_line_prefix": "%m [%r] [%p]",
    "log_min_messages": "INFO",
    "log_statement": "none",
    "maintenance_work_mem": "1GB",
    "max_connections": "600",
    "max_files_per_process": "65535",
    "max_logical_replication_workers": "9",
    "max_prepared_transactions": "200",
    "max_wal_size": "4GB",
    "max_worker_processes": "16",
    "min_wal_size": "512MB",
    "shardman.query_engine_mode": "text",
    "shardman.silk_num_workers": "8",
    "shared_buffers": "4GB",
    "shared_preload_libraries": "postgres_fdw, shardman, biha"
  },
  "strictUserHBA": true,
  "pgHBA": [
    "hostssl all postgres 0.0.0.0/0 cert clientcert=verify-full",
    "hostssl all repluser 0.0.0.0/0 cert clientcert=verify-full",
    "hostssl replication postgres 0.0.0.0/0 cert clientcert=verify-full",
    "hostssl replication postgres ::0/0 cert clientcert=verify-full",
    "hostssl replication repluser 0.0.0.0/0 cert clientcert=verify-full",
    "hostssl replication repluser ::0/0 cert clientcert=verify-full",
    "hostnssl all all 0.0.0.0/0 reject",
    "local postgres postgres scram-sha-256",
    "local replication repluser scram-sha-256"
  ],
  "automaticPgRestart": true,
  "masterDemotionEnabled": false
},
"FDWOptions": {
  "async_capable": "on",
  "batch_size": "100",
  "connect_timeout": "5",
  "fdw_tuple_cost": "0.2",
  "fetch_size": "50000",
  "tcp_user_timeout": "10000"
}
}

```

Note that if the `certVersion` is set to `true`, you cannot modify `pg_hba.conf`.

19.21.4.3. Spec File to Create BiHA Cluster

To create a cluster with topology specified in `sdmspec.json` configuration file, place it under the `Topology` key. The topology format is a key-value pair, with key being a line with the shard name (*shard_name-n*, where *n* is a shard number), and value being the nodes array for the specified shard. Each node must contain the following information:

```

{
  "host": "host",      // name or IP of a node
  "role": "leader",    // node role: leader, follower, referee, referee-with-wal
  "pgPort": 5432,      // PostgreSQL instance port

```

```

    "silkPort": 8000, // silk port
    "bihaPort": 25432, // biha port
    "priority": 10 // node priority, a time interval after which the node can be
suggested as a leader candidate, in milliseconds. The smaller the value, the greater
is priority. Optional parameter.
}

```

Also mind that the specified topology must comply with a number of rules:

- Each shard must contain at least one leader node.
- The number of shard followers must be greater than or equal to the `minSynchronousStandbys` value.
- The `host:pgPort` pair must be unique for each node.
- The ports can't repeat within one node.
- The node priority must be either not specified or be less than 0. The referee or `referee-with-wal` instances must have no priority.

Example:

```

{
  "ConfigVersion": "1",
  "Repfactor": 1,
  "Topology": {
    "shard-1": [
      {
        "host": "host1",
        "role": "leader",
        "pgPort": 5432,
        "silkPort": 8000,
        "bihaPort": 25432,
        "priority": 10
      },
      {
        "host": "host2",
        "role": "follower",
        "pgPort": 5433,
        "silkPort": 8001,
        "bihaPort": 25433,
        "priority": 1000
      },
      {
        "host": "host3",
        "role": "referee",
        "pgPort": 5432,
        "silkPort": 8000,
        "bihaPort": 25432
      }
    ],
    "shard-2": [
      {
        "host": "host2",
        "role": "leader",
        "pgPort": 5432,
        "silkPort": 8000,
        "bihaPort": 25432,
        "priority": 10
      },

```

```

{
  "host": "host1",
  "role": "follower",
  "pgPort": 5433,
  "silkPort": 8001,
  "bihaPort": 25433,
  "priority": 1000
},
{
  "host": "host4",
  "role": "referee-with-wal",
  "pgPort": 5432,
  "silkPort": 8000,
  "bihaPort": 25432
}
],
},
"ShardSpec": {
  "synchronousReplication": true,
  "minSynchronousStandbys": 0,
  ....

```

19.21.4.4. Spec File to Create Custom Logging

To create custom logging, use placeholders from the [shardmanctl config update](#) and [shardmanctl set](#) commands.

Example:

```

{
  "log_directory": "{{dataDir}}",
  "log_filename": "postgres-{{host}}-{{cluster}}-{{shard}}-{{keeperID}}.log",
}

```

19.22. Short Options

For convenience there are also single letter command-line option switches available for some parameters. They are described in [Table 19.4](#). Some of these options exist for historical reasons, and their presence as a single-letter option does not necessarily indicate an endorsement to use the option heavily.

Table 19.4. Short Option Key

Short Option	Equivalent
-B x	shared_buffers = x
-d x	log_min_messages = DEBUG x
-e	datestyle = euro
-fb, -fh, -fi, -fm, -fn, -fo, -fs, -ft	enable_bitmapscan = off , enable_hashjoin = off , enable_indexscan = off, enable_mergejoin = off , enable_nestloop = off, enable_indexonlyscan = off , enable_seqscan = off , enable_tidscan = off
-F	fsync = off
-h x	listen_addresses = x
-i	listen_addresses = '*'
-k x	unix_socket_directories = x
-l	ssl = on

Short Option	Equivalent
-N <i>x</i>	max_connections = <i>x</i>
-O	allow_system_table_mods = on
-p <i>x</i>	port = <i>x</i>
-P	ignore_system_indexes = on
-s	log_statement_stats = on
-S <i>x</i>	work_mem = <i>x</i>
-tpa, -tpl, -te	log_parser_stats = on , log_planner_stats = on , log_executor_stats = on
-W <i>x</i>	post_auth_delay = <i>x</i>

Chapter 20. Client Authentication

When a client application connects to the database server, it specifies which Postgres Pro database user name it wants to connect as, much the same way one logs into a Unix computer as a particular user. Within the SQL environment the active database user name determines access privileges to database objects — see [Chapter 21](#) for more information. Therefore, it is essential to restrict which database users can connect.

Note

As explained in [Chapter 21](#), Postgres Pro actually does privilege management in terms of “roles”. In this chapter, we consistently use *database user* to mean “role with the `LOGIN` privilege”.

Authentication is the process by which the database server establishes the identity of the client, and by extension determines whether the client application (or the user who runs the client application) is permitted to connect with the database user name that was requested.

Postgres Pro offers a number of different client authentication methods. The method used to authenticate a particular client connection can be selected on the basis of (client) host address, database, and user.

Postgres Pro database user names are logically separate from user names of the operating system in which the server runs. If all the users of a particular server also have accounts on the server's machine, it makes sense to assign database user names that match their operating system user names. However, a server that accepts remote connections might have many database users who have no local operating system account, and in such cases there need be no connection between database user names and OS user names.

20.1. The `pg_hba.conf` File

Client authentication is controlled by a configuration file, which traditionally is named `pg_hba.conf` and is stored in the database cluster's data directory. Note that you cannot edit this file manually. (HBA stands for host-based authentication.) A default `pg_hba.conf` file is installed when the data directory is initialized by `initdb`. It is possible to place the authentication configuration file elsewhere, however; see the [hba_file](#) configuration parameter.

The `pg_hba.conf` file is read on start-up and when the main server process receives a `SIGHUP` signal. If you edit the file on an active system, you will need to signal the postmaster (using `pg_ctl reload`, calling the SQL function `pg_reload_conf()`, or using `kill -HUP`) to make it re-read the file.

Note

The preceding statement is not true on Microsoft Windows: there, any changes in the `pg_hba.conf` file are immediately applied by subsequent new connections.

The system view [pg_hba_file_rules](#) can be helpful for pre-testing changes to the `pg_hba.conf` file, or for diagnosing problems if loading of the file did not have the desired effects. Rows in the view with non-null `error` fields indicate problems in the corresponding lines of the file.

The general format of the `pg_hba.conf` file is a set of records, one per line. Blank lines are ignored, as is any text after the `#` comment character. A record can be continued onto the next line by ending the line with a backslash. (Backslashes are not special except at the end of a line.) A record is made up of a number of fields which are separated by spaces and/or tabs. Fields can contain white space if the field value is double-quoted. Quoting one of the keywords in a database, user, or address field (e.g., `all` or `replication`) makes the word lose its special meaning, and just match a database, user, or host with that name. Backslash line continuation applies even within quoted text or comments.

Each authentication record specifies a connection type, a client IP address range (if relevant for the connection type), a database name, a user name, and the authentication method to be used for connections

matching these parameters. The first record with a matching connection type, client address, requested database, and user name is used to perform authentication. There is no “fall-through” or “backup”: if one record is chosen and the authentication fails, subsequent records are not considered. If no record matches, access is denied.

Each record can be an include directive or an authentication record. Include directives specify files that can be included, that contain additional records. The records will be inserted in place of the include directives. Include directives only contain two fields: `include`, `include_if_exists` or `include_dir` directive and the file or directory to be included. The file or directory can be a relative or absolute path, and can be double-quoted. For the `include_dir` form, all files not starting with a `.` and ending with `.conf` will be included. Multiple files within an include directory are processed in file name order (according to C locale rules, i.e., numbers before letters, and uppercase letters before lowercase ones).

A record can have several formats:

<code>local</code>	<code>database</code>	<code>user</code>	<code>auth-method</code>	<code>[auth-options]</code>		
<code>host</code>	<code>database</code>	<code>user</code>	<code>address</code>	<code>auth-method</code>	<code>[auth-options]</code>	
<code>hostssl</code>	<code>database</code>	<code>user</code>	<code>address</code>	<code>auth-method</code>	<code>[auth-options]</code>	
<code>hostnssl</code>	<code>database</code>	<code>user</code>	<code>address</code>	<code>auth-method</code>	<code>[auth-options]</code>	
<code>hostgssenc</code>	<code>database</code>	<code>user</code>	<code>address</code>	<code>auth-method</code>	<code>[auth-options]</code>	
<code>hostnogssenc</code>	<code>database</code>	<code>user</code>	<code>address</code>	<code>auth-method</code>	<code>[auth-options]</code>	
<code>host</code>	<code>database</code>	<code>user</code>	<code>IP-address</code>	<code>IP-mask</code>	<code>auth-method</code>	<code>[auth-options]</code>
<code>hostssl</code>	<code>database</code>	<code>user</code>	<code>IP-address</code>	<code>IP-mask</code>	<code>auth-method</code>	<code>[auth-options]</code>
<code>hostnssl</code>	<code>database</code>	<code>user</code>	<code>IP-address</code>	<code>IP-mask</code>	<code>auth-method</code>	<code>[auth-options]</code>
<code>hostgssenc</code>	<code>database</code>	<code>user</code>	<code>IP-address</code>	<code>IP-mask</code>	<code>auth-method</code>	<code>[auth-options]</code>
<code>hostnogssenc</code>	<code>database</code>	<code>user</code>	<code>IP-address</code>	<code>IP-mask</code>	<code>auth-method</code>	<code>[auth-options]</code>
<code>include</code>	<code>file</code>					
<code>include_if_exists</code>	<code>file</code>					
<code>include_dir</code>	<code>directory</code>					

The meaning of the fields is as follows:

`local`

This record matches connection attempts using Unix-domain sockets. Without a record of this type, Unix-domain socket connections are disallowed.

`host`

This record matches connection attempts made using TCP/IP. `host` records match SSL or non-SSL connection attempts as well as GSSAPI encrypted or non-GSSAPI encrypted connection attempts.

Note

Remote TCP/IP connections will not be possible unless the server is started with an appropriate value for the `listen_addresses` configuration parameter, since the default behavior is to listen for TCP/IP connections only on the local loopback address `localhost`.

`hostssl`

This record matches connection attempts made using TCP/IP, but only when the connection is made with SSL encryption.

To make use of this option the server must be built with SSL support. Furthermore, SSL must be enabled by setting the `ssl` configuration parameter (see [Section 18.9](#) for more information). Otherwise, the `hostssl` record is ignored except for logging a warning that it cannot match any connections.

hostnossl

This record type has the opposite behavior of *hostssl*; it only matches connection attempts made over TCP/IP that do not use SSL.

hostgssenc

This record matches connection attempts made using TCP/IP, but only when the connection is made with GSSAPI encryption.

To make use of this option the server must be built with GSSAPI support. Otherwise, the *hostgssenc* record is ignored except for logging a warning that it cannot match any connections.

hostnogssenc

This record type has the opposite behavior of *hostgssenc*; it only matches connection attempts made over TCP/IP that do not use GSSAPI encryption.

database

Specifies which database name(s) this record matches. The value *all* specifies that it matches all databases. The value *sameuser* specifies that the record matches if the requested database has the same name as the requested user. The value *samerole* specifies that the requested user must be a member of the role with the same name as the requested database. (*samegroup* is an obsolete but still accepted spelling of *samerole*.) Superusers are not considered to be members of a role for the purposes of *samerole* unless they are explicitly members of the role, directly or indirectly, and not just by virtue of being a superuser. The value *replication* specifies that the record matches if a physical replication connection is requested, however, it doesn't match with logical replication connections. Note that physical replication connections do not specify any particular database whereas logical replication connections do specify it. Otherwise, this is the name of a specific Postgres Pro database or a regular expression. Multiple database names and/or regular expressions can be supplied by separating them with commas.

If the database name starts with a slash (/), the remainder of the name is treated as a regular expression. (See [Section 9.7.3.1](#) for details of Postgres Pro's regular expression syntax.)

A separate file containing database names and/or regular expressions can be specified by preceding the file name with @.

user

Specifies which database user name(s) this record matches. The value *all* specifies that it matches all users. Otherwise, this is either the name of a specific database user, a regular expression (when starting with a slash (/), or a group name preceded by +. (Recall that there is no real distinction between users and groups in Postgres Pro; a + mark really means “match any of the roles that are directly or indirectly members of this role”, while a name without a + mark matches only that specific role.) For this purpose, a superuser is only considered to be a member of a role if they are explicitly a member of the role, directly or indirectly, and not just by virtue of being a superuser. Multiple user names and/or regular expressions can be supplied by separating them with commas.

If the user name starts with a slash (/), the remainder of the name is treated as a regular expression. (See [Section 9.7.3.1](#) for details of Postgres Pro's regular expression syntax.)

A separate file containing user names and/or regular expressions can be specified by preceding the file name with @.

address

Specifies the client machine address(es) that this record matches. This field can contain either a host name, an IP address range, or one of the special key words mentioned below.

An IP address range is specified using standard numeric notation for the range's starting address, then a slash (/) and a CIDR mask length. The mask length indicates the number of high-order bits of

the client IP address that must match. Bits to the right of this should be zero in the given IP address. There must not be any white space between the IP address, the /, and the CIDR mask length.

Typical examples of an IPv4 address range specified this way are `172.20.143.89/32` for a single host, or `172.20.143.0/24` for a small network, or `10.6.0.0/16` for a larger one. An IPv6 address range might look like `::1/128` for a single host (in this case the IPv6 loopback address) or `fe80::7a31:c1f-f:0000:0000/96` for a small network. `0.0.0.0/0` represents all IPv4 addresses, and `:::0/0` represents all IPv6 addresses. To specify a single host, use a mask length of 32 for IPv4 or 128 for IPv6. In a network address, do not omit trailing zeroes.

An entry given in IPv4 format will match only IPv4 connections, and an entry given in IPv6 format will match only IPv6 connections, even if the represented address is in the IPv4-in-IPv6 range.

You can also write `all` to match any IP address, `samehost` to match any of the server's own IP addresses, or `sameonet` to match any address in any subnet that the server is directly connected to.

If a host name is specified (anything that is not an IP address range or a special key word is treated as a host name), that name is compared with the result of a reverse name resolution of the client's IP address (e.g., reverse DNS lookup, if DNS is used). Host name comparisons are case insensitive. If there is a match, then a forward name resolution (e.g., forward DNS lookup) is performed on the host name to check whether any of the addresses it resolves to are equal to the client's IP address. If both directions match, then the entry is considered to match. (The host name that is used in `pg_hba.conf` should be the one that address-to-name resolution of the client's IP address returns, otherwise the line won't be matched. Some host name databases allow associating an IP address with multiple host names, but the operating system will only return one host name when asked to resolve an IP address.)

A host name specification that starts with a dot (.) matches a suffix of the actual host name. So `.example.com` would match `foo.example.com` (but not just `example.com`).

When host names are specified in `pg_hba.conf`, you should make sure that name resolution is reasonably fast. It can be of advantage to set up a local name resolution cache such as `nsd`. Also, you may wish to enable the configuration parameter `log_hostname` to see the client's host name instead of the IP address in the log.

These fields do not apply to `local` records.

Note

Users sometimes wonder why host names are handled in this seemingly complicated way, with two name resolutions including a reverse lookup of the client's IP address. This complicates use of the feature in case the client's reverse DNS entry is not set up or yields some undesirable host name. It is done primarily for efficiency: this way, a connection attempt requires at most two resolver lookups, one reverse and one forward. If there is a resolver problem with some address, it becomes only that client's problem. A hypothetical alternative implementation that only did forward lookups would have to resolve every host name mentioned in `pg_hba.conf` during every connection attempt. That could be quite slow if many names are listed. And if there is a resolver problem with one of the host names, it becomes everyone's problem.

Also, a reverse lookup is necessary to implement the suffix matching feature, because the actual client host name needs to be known in order to match it against the pattern.

Note that this behavior is consistent with other popular implementations of host name-based access control, such as the Apache HTTP Server and TCP Wrappers.

IP-address

IP-mask

These two fields can be used as an alternative to the `IP-address/mask-length` notation. Instead of specifying the mask length, the actual mask is specified in a separate column. For example,

255.0.0.0 represents an IPv4 CIDR mask length of 8, and 255.255.255.255 represents a CIDR mask length of 32.

These fields do not apply to `local` records.

auth-method

Specifies the authentication method to use when a connection matches this record. The possible choices are summarized here; details are in [Section 20.3](#). All the options are lower case and treated case sensitively, so even acronyms like `ldap` must be specified as lower case.

trust

Allow the connection unconditionally. This method allows anyone that can connect to the Postgres Pro database server to login as any Postgres Pro user they wish, without the need for a password or any other authentication. See [Section 20.4](#) for details.

reject

Reject the connection unconditionally. This is useful for “filtering out” certain hosts from a group, for example a `reject` line could block a specific host from connecting, while a later line allows the remaining hosts in a specific network to connect.

scram-sha-256

Perform SCRAM-SHA-256 authentication to verify the user's password. See [Section 20.5](#) for details.

md5

Perform SCRAM-SHA-256 or MD5 authentication to verify the user's password. See [Section 20.5](#) for details.

password

Require the client to supply an unencrypted password for authentication. Since the password is sent in clear text over the network, this should not be used on untrusted networks. See [Section 20.5](#) for details.

gss

Use GSSAPI to authenticate the user. This is only available for TCP/IP connections. See [Section 20.6](#) for details. It can be used in conjunction with GSSAPI encryption.

sspi

Use SSPI to authenticate the user. This is only available on Windows. See [Section 20.7](#) for details.

ident

Obtain the operating system user name of the client by contacting the ident server on the client and check if it matches the requested database user name. Ident authentication can only be used on TCP/IP connections. When specified for local connections, peer authentication will be used instead. See [Section 20.8](#) for details.

peer

Obtain the client's operating system user name from the operating system and check if it matches the requested database user name. This is only available for local connections. See [Section 20.9](#) for details.

ldap

Authenticate using an LDAP server. See [Section 20.10](#) for details.

radius

Authenticate using a RADIUS server. See [Section 20.11](#) for details.

`cert`

Authenticate using SSL client certificates. See [Section 20.12](#) for details.

`pam`

Authenticate using the Pluggable Authentication Modules (PAM) service provided by the operating system. See [Section 20.13](#) for details.

`bsd`

Authenticate using the BSD Authentication service provided by the operating system. See [Section 20.14](#) for details.

auth-options

After the *auth-method* field, there can be field(s) of the form *name=value* that specify options for the authentication method. Details about which options are available for which authentication methods appear below.

In addition to the method-specific options listed below, there is a method-independent authentication option `clientcert`, which can be specified in any `hostssl` record. This option can be set to `verify-ca` or `verify-full`. Both options require the client to present a valid (trusted) SSL certificate, while `verify-full` additionally enforces that the `cn` (Common Name) in the certificate matches the username or an applicable mapping. This behavior is similar to the `cert` authentication method (see [Section 20.12](#)) but enables pairing the verification of client certificates with any authentication method that supports `hostssl` entries.

On any record using client certificate authentication (i.e. one using the `cert` authentication method or one using the `clientcert` option), you can specify which part of the client certificate credentials to match using the `clientname` option. This option can have one of two values. If you specify `clientname=CN`, which is the default, the username is matched against the certificate's Common Name (CN). If instead you specify `clientname=DN` the username is matched against the entire Distinguished Name (DN) of the certificate. This option is probably best used in conjunction with a username map. The comparison is done with the DN in [RFC 2253](#) format. To see the DN of a client certificate in this format, do

```
openssl x509 -in myclient.crt -noout -subject -nameopt RFC2253 | sed "s/^subject=/"
```

Care needs to be taken when using this option, especially when using regular expression matching against the DN.

`include`

This line will be replaced by the contents of the given file.

`include_if_exists`

This line will be replaced by the content of the given file if the file exists. Otherwise, a message is logged to indicate that the file has been skipped.

`include_dir`

This line will be replaced by the contents of all the files found in the directory, if they don't start with a `.` and end with `.conf`, processed in file name order (according to C locale rules, i.e., numbers before letters, and uppercase letters before lowercase ones).

Files included by `@` constructs are read as lists of names, which can be separated by either whitespace or commas. Comments are introduced by `#`, just as in `pg_hba.conf`, and nested `@` constructs are allowed. Unless the file name following `@` is an absolute path, it is taken to be relative to the directory containing the referencing file.

Since the `pg_hba.conf` records are examined sequentially for each connection attempt, the order of the records is significant. Typically, earlier records will have tight connection match parameters and weaker authentication methods, while later records will have looser match parameters and stronger authentication methods. For example, one might wish to use `trust` authentication for local TCP/IP connections but

require a password for remote TCP/IP connections. In this case a record specifying `trust` authentication for connections from 127.0.0.1 would appear before a record specifying password authentication for a wider range of allowed client IP addresses.

Tip

To connect to a particular database, a user must not only pass the `pg_hba.conf` checks, but must have the `CONNECT` privilege for the database. If you wish to restrict which users can connect to which databases, it's usually easier to control this by granting/revoking `CONNECT` privilege than to put the rules in `pg_hba.conf` entries.

Some examples of `pg_hba.conf` entries are shown in [Example 20.1](#). See the next section for details on the different authentication methods.

Example 20.1. Example `pg_hba.conf` Entries

```
# Allow any user on the local system to connect to any database with
# any database user name using Unix-domain sockets (the default for local
# connections).
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
local      all             all
trust

# The same using local loopback TCP/IP connections.
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
host       all             all        127.0.0.1/32  trust

# The same as the previous line, but using a separate netmask column
#
# TYPE      DATABASE      USER      IP-ADDRESS    IP-MASK      METHOD
host       all             all        127.0.0.1     255.255.255.255  trust

# The same over IPv6.
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
host       all             all        ::1/128      trust

# The same using a host name (would typically cover both IPv4 and IPv6).
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
host       all             all        localhost    trust

# The same using a regular expression for DATABASE, that allows connection
# to any databases with a name beginning with "db" and finishing with a
# number using two to four digits (like "db1234" or "db12").
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
host       "^db\d{2,4}$"  all        localhost    trust

# Allow any user from any host with IP address 192.168.93.x to connect
# to database "postgres" as the same user name that ident reports for
# the connection (typically the operating system user name).
#
# TYPE      DATABASE      USER      ADDRESS      METHOD
host       postgres     all        192.168.93.0/24  ident

# Allow any user from host 192.168.12.10 to connect to database
```

```
# "postgres" if the user's password is correctly supplied.
#
# TYPE  DATABASE          USER                ADDRESS              METHOD
host    postgres          all                 192.168.12.10/32    scram-sha-256

# Allow any user from hosts in the example.com domain to connect to
# any database if the user's password is correctly supplied.
#
# Require SCRAM authentication for most users, but make an exception
# for user 'mike', who uses an older client that doesn't support SCRAM
# authentication.
#
# TYPE  DATABASE          USER                ADDRESS              METHOD
host    all               mike                .example.com         md5
host    all               all                 .example.com         scram-sha-256

# In the absence of preceding "host" lines, these three lines will
# reject all connections from 192.168.54.1 (since that entry will be
# matched first), but allow GSSAPI-encrypted connections from anywhere else
# on the Internet. The zero mask causes no bits of the host IP address to
# be considered, so it matches any host. Unencrypted GSSAPI connections
# (which "fall through" to the third line since "hostgssenc" only matches
# encrypted GSSAPI connections) are allowed, but only from 192.168.12.10.
#
# TYPE  DATABASE          USER                ADDRESS              METHOD
host    all               all                 192.168.54.1/32     reject
hostgssenc all           all                 0.0.0.0/0            gss
host    all               all                 192.168.12.10/32    gss

# Allow users from 192.168.x.x hosts to connect to any database, if
# they pass the ident check. If, for example, ident says the user is
# "bryanh" and he requests to connect as Postgres Pro user "guest1", the
# connection is allowed if there is an entry in pg_ident.conf for map
# "omicron" that says "bryanh" is allowed to connect as "guest1".
#
# TYPE  DATABASE          USER                ADDRESS              METHOD
host    all               all                 192.168.0.0/16       ident map=omicron

# If these are the only four lines for local connections, they will
# allow local users to connect only to their own databases (databases
# with the same name as their database user name) except for users whose
# name end with "helpdesk", administrators and members of role "support",
# who can connect to all databases. The file $PGDATA/admins contains a
# list of names of administrators. Passwords are required in all cases.
#
# TYPE  DATABASE          USER                ADDRESS              METHOD
local   sameuser        all                 all                  md5
local   all              /^.*helpdesk$      all                  md5
local   all              @admins             all                  md5
local   all              +support            all                  md5

# The last two lines above can be combined into a single line:
local   all              @admins,+support    all                  md5

# The database column can also use lists and file names:
local   db1,db2,@demodbs all                  md5
```


20.2. User Name Maps

When using an external authentication system such as Ident or GSSAPI, the name of the operating system user that initiated the connection might not be the same as the database user (role) that is to be used. In this case, a user name map can be applied to map the operating system user name to a database user. To use user name mapping, specify `map=map-name` in the options field in `pg_hba.conf`. This option is supported for all authentication methods that receive external user names. Since different mappings might be needed for different connections, the name of the map to be used is specified in the `map-name` parameter in `pg_hba.conf` to indicate which map to use for each individual connection.

User name maps are defined in the ident map file, which by default is named `pg_ident.conf` and is stored in the cluster's data directory. (It is possible to place the map file elsewhere, however; see the [ident_file](#) configuration parameter.) The ident map file contains lines of the general forms:

```
map-name system-username database-username
include file
include_if_exists file
include_dir directory
```

Comments, whitespace and line continuations are handled in the same way as in `pg_hba.conf`. The `map-name` is an arbitrary name that will be used to refer to this mapping in `pg_hba.conf`. The other two fields specify an operating system user name and a matching database user name. The same `map-name` can be used repeatedly to specify multiple user-mappings within a single map.

As for `pg_hba.conf`, the lines in this file can be include directives, following the same rules.

The `pg_ident.conf` file is read on start-up and when the main server process receives a SIGHUP signal. If you edit the file on an active system, you will need to signal the postmaster (using `pg_ctl reload`, calling the SQL function `pg_reload_conf()`, or using `kill -HUP`) to make it re-read the file.

The system view `pg_ident_file_mappings` can be helpful for pre-testing changes to the `pg_ident.conf` file, or for diagnosing problems if loading of the file did not have the desired effects. Rows in the view with non-null `error` fields indicate problems in the corresponding lines of the file.

There is no restriction regarding how many database users a given operating system user can correspond to, nor vice versa. Thus, entries in a map should be thought of as meaning “this operating system user is allowed to connect as this database user”, rather than implying that they are equivalent. The connection will be allowed if there is any map entry that pairs the user name obtained from the external authentication system with the database user name that the user has requested to connect as. The value `all` can be used as the `database-username` to specify that if the `system-username` matches, then this user is allowed to log in as any of the existing database users. Quoting `all` makes the keyword lose its special meaning.

If the `database-username` begins with a `+` character, then the operating system user can login as any user belonging to that role, similarly to how user names beginning with `+` are treated in `pg_hba.conf`. Thus, a `+` mark means “match any of the roles that are directly or indirectly members of this role”, while a name without a `+` mark matches only that specific role. Quoting a username starting with a `+` makes the `+` lose its special meaning.

If the `system-username` field starts with a slash (`/`), the remainder of the field is treated as a regular expression. (See [Section 9.7.3.1](#) for details of Postgres Pro's regular expression syntax.) The regular expression can include a single capture, or parenthesized subexpression, which can then be referenced in the `database-username` field as `\1` (backslash-one). This allows the mapping of multiple user names in a single line, which is particularly useful for simple syntax substitutions. For example, these entries

```
mymap    /^(.*)@mydomain\.com$      \1
mymap    /^(.*)@otherdomain\.com$   guest
```

will remove the domain part for users with system user names that end with `@mydomain.com`, and allow any user whose system name ends with `@otherdomain.com` to log in as `guest`. Quoting a `database-username` containing `\1` *does not* make `\1` lose its special meaning.

If the `database-username` field starts with a slash (/), the remainder of the field is treated as a regular expression (see [Section 9.7.3.1](#) for details of Postgres Pro's regular expression syntax). It is not possible to use `\1` to use a capture from regular expression on `system-username` for a regular expression on `database-username`.

Tip

Keep in mind that by default, a regular expression can match just part of a string. It's usually wise to use `^` and `$`, as shown in the above example, to force the match to be to the entire system user name.

A `pg_ident.conf` file that could be used in conjunction with the `pg_hba.conf` file in [Example 20.1](#) is shown in [Example 20.2](#). In this example, anyone logged in to a machine on the 192.168 network that does not have the operating system user name `bryanh`, `ann`, or `robert` would not be granted access. Unix user `robert` would only be allowed access when he tries to connect as Postgres Pro user `bob`, not as `robert` or anyone else. `ann` would only be allowed to connect as `ann`. User `bryanh` would be allowed to connect as either `bryanh` or as `guest1`.

Example 20.2. An Example `pg_ident.conf` File

```
# MAPNAME          SYSTEM-USERNAME      PG-USERNAME

omicron            bryanh                bryanh
omicron            ann                  ann
# bob has user name robert on these machines
omicron            robert                bob
# bryanh can also connect as guest1
omicron            bryanh                guest1
```

20.3. Authentication Methods

Postgres Pro provides various methods for authenticating users:

- [Trust authentication](#), which simply trusts that users are who they say they are.
- [Password authentication](#), which requires that users send a password.
- [GSSAPI authentication](#), which relies on a GSSAPI-compatible security library. Typically this is used to access an authentication server such as a Kerberos or Microsoft Active Directory server.
- [SSPI authentication](#), which uses a Windows-specific protocol similar to GSSAPI.
- [Ident authentication](#), which relies on an “Identification Protocol” ([RFC 1413](#)) service on the client's machine. (On local Unix-socket connections, this is treated as peer authentication.)
- [Peer authentication](#), which relies on operating system facilities to identify the process at the other end of a local connection. This is not supported for remote connections.
- [LDAP authentication](#), which relies on an LDAP authentication server.
- [RADIUS authentication](#), which relies on a RADIUS authentication server.
- [Certificate authentication](#), which requires an SSL connection and authenticates users by checking the SSL certificate they send.
- [PAM authentication](#), which relies on a PAM (Pluggable Authentication Modules) library.
- [BSD authentication](#), which relies on the BSD Authentication framework (currently available only on OpenBSD).

Peer authentication is usually recommendable for local connections, though trust authentication might be sufficient in some circumstances. Password authentication is the easiest choice for remote connections. All the other options require some kind of external security infrastructure (usually an authentication server or a certificate authority for issuing SSL certificates), or are platform-specific.

The following sections describe each of these authentication methods in more detail.

20.4. Trust Authentication

When `trust` authentication is specified, Postgres Pro assumes that anyone who can connect to the server is authorized to access the database with whatever database user name they specify (even superuser names). Of course, restrictions made in the `database` and `user` columns still apply. This method should only be used when there is adequate operating-system-level protection on connections to the server.

`trust` authentication is appropriate and very convenient for local connections on a single-user workstation. It is usually *not* appropriate by itself on a multiuser machine. However, you might be able to use `trust` even on a multiuser machine, if you restrict access to the server's Unix-domain socket file using file-system permissions. To do this, set the `unix_socket_permissions` (and possibly `unix_socket_group`) configuration parameters as described in [Section 19.3](#). Or you could set the `unix_socket_directories` configuration parameter to place the socket file in a suitably restricted directory.

Setting file-system permissions only helps for Unix-socket connections. Local TCP/IP connections are not restricted by file-system permissions. Therefore, if you want to use file-system permissions for local security, remove the `host ... 127.0.0.1 ...` line from `pg_hba.conf`, or change it to a non-`trust` authentication method.

`trust` authentication is only suitable for TCP/IP connections if you trust every user on every machine that is allowed to connect to the server by the `pg_hba.conf` lines that specify `trust`. It is seldom reasonable to use `trust` for any TCP/IP connections other than those from localhost (127.0.0.1).

20.5. Password Authentication

There are several password-based authentication methods. These methods operate similarly but differ in how the users' passwords are stored on the server and how the password provided by a client is sent across the connection.

`scram-sha-256`

The method `scram-sha-256` performs SCRAM-SHA-256 authentication, as described in [RFC 7677](#). It is a challenge-response scheme that prevents password sniffing on untrusted connections and supports storing passwords on the server in a cryptographically hashed form that is thought to be secure.

This is the most secure of the currently provided methods, but it is not supported by older client libraries.

`md5`

The method `md5` uses a custom less secure challenge-response mechanism. It prevents password sniffing and avoids storing passwords on the server in plain text but provides no protection if an attacker manages to steal the password hash from the server. Also, the MD5 hash algorithm is nowadays no longer considered secure against determined attacks.

To ease transition from the `md5` method to the newer SCRAM method, if `md5` is specified as a method in `pg_hba.conf` but the user's password on the server is encrypted for SCRAM (see below), then SCRAM-based authentication will automatically be chosen instead.

`password`

The method `password` sends the password in clear-text and is therefore vulnerable to password “sniffing” attacks. It should always be avoided if possible. If the connection is protected by SSL encryption then `password` can be used safely, though. (Though SSL certificate authentication might be a better choice if one is depending on using SSL).

Postgres Pro database passwords are separate from operating system user passwords. The password for each database user is stored in the `pg_authid` system catalog. Passwords can be managed with the SQL commands [CREATE ROLE](#) and [ALTER ROLE](#), e.g., `CREATE ROLE foo WITH LOGIN PASSWORD 'secret'`,

or the `psql` command `\password`. If no password has been set up for a user, the stored password is null and password authentication will always fail for that user.

The availability of the different password-based authentication methods depends on how a user's password on the server is encrypted (or hashed, more accurately). This is controlled by the configuration parameter `password_encryption` at the time the password is set. If a password was encrypted using the `scram-sha-256` setting, then it can be used for the authentication methods `scram-sha-256` and `password` (but password transmission will be in plain text in the latter case). The authentication method specification `md5` will automatically switch to using the `scram-sha-256` method in this case, as explained above, so it will also work. If a password was encrypted using the `md5` setting, then it can be used only for the `md5` and `password` authentication method specifications (again, with the password transmitted in plain text in the latter case). (Previous Postgres Pro releases supported storing the password on the server in plain text. This is no longer possible.) To check the currently stored password hashes, see the system catalog `pg_authid`.

To upgrade an existing installation from `md5` to `scram-sha-256`, after having ensured that all client libraries in use are new enough to support SCRAM, set `password_encryption = 'scram-sha-256'` in `postgresql.conf`, make all users set new passwords, and change the authentication method specifications in `pg_hba.conf` to `scram-sha-256`.

20.6. GSSAPI Authentication

GSSAPI is an industry-standard protocol for secure authentication defined in [RFC 2743](#). Postgres Pro supports GSSAPI for authentication, communications encryption, or both. GSSAPI provides automatic authentication (single sign-on) for systems that support it. The authentication itself is secure. If GSSAPI encryption or SSL encryption is used, the data sent along the database connection will be encrypted; otherwise, it will not.

GSSAPI support has to be enabled when Postgres Pro is built.

When GSSAPI uses Kerberos, it uses a standard service principal (authentication identity) name in the format `servicename/hostname@realm`. The principal name used by a particular installation is not encoded in the Postgres Pro server in any way; rather it is specified in the `keytab` file that the server reads to determine its identity. If multiple principals are listed in the `keytab` file, the server will accept any one of them. The server's realm name is the preferred realm specified in the Kerberos configuration file(s) accessible to the server.

When connecting, the client must know the principal name of the server it intends to connect to. The `servicename` part of the principal is ordinarily `postgres`, but another value can be selected via libpq's `krb_srvname` connection parameter. The `hostname` part is the fully qualified host name that libpq is told to connect to. The realm name is the preferred realm specified in the Kerberos configuration file(s) accessible to the client.

The client will also have a principal name for its own identity (and it must have a valid ticket for this principal). To use GSSAPI for authentication, the client principal must be associated with a Postgres Pro database user name. The `pg_ident.conf` configuration file can be used to map principals to user names; for example, `pgusername@realm` could be mapped to just `pgusername`. Alternatively, you can use the full `username@realm` principal as the role name in Postgres Pro without any mapping.

Postgres Pro also supports mapping client principals to user names by just stripping the realm from the principal. This method is supported for backwards compatibility and is strongly discouraged as it is then impossible to distinguish different users with the same user name but coming from different realms. To enable this, set `include_realm` to 0. For simple single-realm installations, doing that combined with setting the `krb_realm` parameter (which checks that the principal's realm matches exactly what is in the `krb_realm` parameter) is still secure; but this is a less capable approach compared to specifying an explicit mapping in `pg_ident.conf`.

The location of the server's `keytab` file is specified by the `krb_server_keyfile` configuration parameter. For security reasons, it is recommended to use a separate `keytab` just for the Postgres Pro server rather

than allowing the server to read the system keytab file. Make sure that your server keytab file is readable (and preferably only readable, not writable) by the Postgres Pro server account. (See also [Section 18.1.](#))

The keytab file is generated using the Kerberos software; see the Kerberos documentation for details. The following example shows doing this using the `kadmin` tool of MIT Kerberos:

```
kadmin% addprinc -randkey postgres/server.my.domain.org
kadmin% ktadd -k krb5.keytab postgres/server.my.domain.org
```

The following authentication options are supported for the GSSAPI authentication method:

`include_realm`

If set to 0, the realm name from the authenticated user principal is stripped off before being passed through the user name mapping ([Section 20.2](#)). This is discouraged and is primarily available for backwards compatibility, as it is not secure in multi-realm environments unless `krb_realm` is also used. It is recommended to leave `include_realm` set to the default (1) and to provide an explicit mapping in `pg_ident.conf` to convert principal names to Postgres Pro user names.

`map`

Allows mapping from client principals to database user names. See [Section 20.2](#) for details. For a GSSAPI/Kerberos principal, such as `username@EXAMPLE.COM` (or, less commonly, `username/hostbased@EXAMPLE.COM`), the user name used for mapping is `username@EXAMPLE.COM` (or `username/hostbased@EXAMPLE.COM`, respectively), unless `include_realm` has been set to 0, in which case `username` (or `username/hostbased`) is what is seen as the system user name when mapping.

`krb_realm`

Sets the realm to match user principal names against. If this parameter is set, only users of that realm will be accepted. If it is not set, users of any realm can connect, subject to whatever user name mapping is done.

In addition to these settings, which can be different for different `pg_hba.conf` entries, there is the server-wide `krb_caseins_users` configuration parameter. If that is set to true, client principals are matched to user map entries case-insensitively. `krb_realm`, if set, is also matched case-insensitively.

20.7. SSPI Authentication

SSPI is a Windows technology for secure authentication with single sign-on. Postgres Pro will use SSPI in `negotiate` mode, which will use Kerberos when possible and automatically fall back to NTLM in other cases. SSPI and GSSAPI interoperate as clients and servers, e.g., an SSPI client can authenticate to an GSSAPI server. It is recommended to use SSPI on Windows clients and servers and GSSAPI on non-Windows platforms.

When using Kerberos authentication, SSPI works the same way GSSAPI does; see [Section 20.6](#) for details.

The following configuration options are supported for SSPI:

`include_realm`

If set to 0, the realm name from the authenticated user principal is stripped off before being passed through the user name mapping ([Section 20.2](#)). This is discouraged and is primarily available for backwards compatibility, as it is not secure in multi-realm environments unless `krb_realm` is also used. It is recommended to leave `include_realm` set to the default (1) and to provide an explicit mapping in `pg_ident.conf` to convert principal names to Postgres Pro user names.

`compat_realm`

If set to 1, the domain's SAM-compatible name (also known as the NetBIOS name) is used for the `include_realm` option. This is the default. If set to 0, the true realm name from the Kerberos user principal name is used.

Do not disable this option unless your server runs under a domain account (this includes virtual service accounts on a domain member system) and all clients authenticating through SSPI are also using domain accounts, or authentication will fail.

`upn_username`

If this option is enabled along with `compat_realm`, the user name from the Kerberos UPN is used for authentication. If it is disabled (the default), the SAM-compatible user name is used. By default, these two names are identical for new user accounts.

Note that `libpq` uses the SAM-compatible name if no explicit user name is specified. If you use `libpq` or a driver based on it, you should leave this option disabled or explicitly specify user name in the connection string.

`map`

Allows for mapping between system and database user names. See [Section 20.2](#) for details. For an SSPI/Kerberos principal, such as `username@EXAMPLE.COM` (or, less commonly, `username/host-based@EXAMPLE.COM`), the user name used for mapping is `username@EXAMPLE.COM` (or `username/hostbased@EXAMPLE.COM`, respectively), unless `include_realm` has been set to 0, in which case `username` (or `username/hostbased`) is what is seen as the system user name when mapping.

`krb_realm`

Sets the realm to match user principal names against. If this parameter is set, only users of that realm will be accepted. If it is not set, users of any realm can connect, subject to whatever user name mapping is done.

20.8. Ident Authentication

The ident authentication method works by obtaining the client's operating system user name from an ident server and using it as the allowed database user name (with an optional user name mapping). This is only supported on TCP/IP connections.

Note

When ident is specified for a local (non-TCP/IP) connection, peer authentication (see [Section 20.9](#)) will be used instead.

The following configuration options are supported for `ident`:

`map`

Allows for mapping between system and database user names. See [Section 20.2](#) for details.

The “Identification Protocol” is described in [RFC 1413](#). Virtually every Unix-like operating system ships with an ident server that listens on TCP port 113 by default. The basic functionality of an ident server is to answer questions like “What user initiated the connection that goes out of your port *x* and connects to my port *y*?”. Since Postgres Pro knows both *x* and *y* when a physical connection is established, it can interrogate the ident server on the host of the connecting client and can theoretically determine the operating system user for any given connection.

The drawback of this procedure is that it depends on the integrity of the client: if the client machine is untrusted or compromised, an attacker could run just about any program on port 113 and return any user name they choose. This authentication method is therefore only appropriate for closed networks where each client machine is under tight control and where the database and system administrators operate in close contact. In other words, you must trust the machine running the ident server. Heed the warning:

The Identification Protocol is not intended as an authorization or access control protocol.

Some ident servers have a nonstandard option that causes the returned user name to be encrypted, using a key that only the originating machine's administrator knows. This option *must not* be used when using the ident server with Postgres Pro, since Postgres Pro does not have any way to decrypt the returned string to determine the actual user name.

20.9. Peer Authentication

The peer authentication method works by obtaining the client's operating system user name from the kernel and using it as the allowed database user name (with optional user name mapping). This method is only supported on local connections.

The following configuration options are supported for `peer`:

`map`

Allows for mapping between system and database user names. See [Section 20.2](#) for details.

Peer authentication is only available on operating systems providing the `getpeereid()` function, the `SO_PEERCREC` socket parameter, or similar mechanisms. Currently that includes Linux, most flavors of BSD including macOS, and Solaris.

20.10. LDAP Authentication

This authentication method operates similarly to `password` except that it uses LDAP as the password verification method. LDAP is used only to validate the user name/password pairs. Therefore the user must already exist in the database before LDAP can be used for authentication.

LDAP authentication can operate in two modes. In the first mode, which we will call the simple bind mode, the server will bind to the distinguished name constructed as `prefix username suffix`. Typically, the `prefix` parameter is used to specify `cn=`, or `DOMAIN\` in an Active Directory environment. `suffix` is used to specify the remaining part of the DN in a non-Active Directory environment.

In the second mode, which we will call the search+bind mode, the server first binds to the LDAP directory with a fixed user name and password, specified with `ldapbinddn` and `ldapbindpasswd`, and performs a search for the user trying to log in to the database. If no user and password is configured, an anonymous bind will be attempted to the directory. The search will be performed over the subtree at `ldapbasedn`, and will try to do an exact match of the attribute specified in `ldapsearchattribute`. Once the user has been found in this search, the server re-binds to the directory as this user, using the password specified by the client, to verify that the login is correct. This mode is the same as that used by LDAP authentication schemes in other software, such as Apache `mod_authnz_ldap` and `pam_ldap`. This method allows for significantly more flexibility in where the user objects are located in the directory, but will cause two additional requests to the LDAP server to be made.

The following configuration options are used in both modes:

`ldapserver`

Names or IP addresses of LDAP servers to connect to. Multiple servers may be specified, separated by spaces.

`ldapport`

Port number on LDAP server to connect to. If no port is specified, the LDAP library's default port setting will be used.

`ldapscheme`

Set to `ldaps` to use LDAPS. This is a non-standard way of using LDAP over SSL, supported by some LDAP server implementations. See also the `ldaptls` option for an alternative.

ldaptls

Set to 1 to make the connection between Postgres Pro and the LDAP server use TLS encryption. This uses the `StartTLS` operation per [RFC 4513](#). See also the `ldapscheme` option for an alternative.

Note that using `ldapscheme` or `ldaptls` only encrypts the traffic between the Postgres Pro server and the LDAP server. The connection between the Postgres Pro server and the Postgres Pro client will still be unencrypted unless SSL is used there as well.

The following options are used in simple bind mode only:

ldapprefix

String to prepend to the user name when forming the DN to bind as, when doing simple bind authentication.

ldapsuffix

String to append to the user name when forming the DN to bind as, when doing simple bind authentication.

The following options are used in search+bind mode only:

ldapbasedn

Root DN to begin the search for the user in, when doing search+bind authentication.

ldapbinddn

DN of user to bind to the directory with to perform the search when doing search+bind authentication.

ldapbindpasswd

Password for user to bind to the directory with to perform the search when doing search+bind authentication.

ldapsearchattribute

Attribute to match against the user name in the search when doing search+bind authentication. If no attribute is specified, the `uid` attribute will be used.

ldapsearchfilter

The search filter to use when doing search+bind authentication. Occurrences of `$username` will be replaced with the user name. This allows for more flexible search filters than `ldapsearchattribute`.

ldapurl

An [RFC 4516](#) LDAP URL. This is an alternative way to write some of the other LDAP options in a more compact and standard form. The format is

```
ldap[s]://host[:port]/basedn[?[attribute][?[scope][?[filter]]]
```

scope must be one of `base`, `one`, `sub`, typically the last. (The default is `base`, which is normally not useful in this application.) *attribute* can nominate a single attribute, in which case it is used as a value for `ldapsearchattribute`. If *attribute* is empty then *filter* can be used as a value for `ldapsearchfilter`.

The URL scheme `ldaps` chooses the LDAPS method for making LDAP connections over SSL, equivalent to using `ldapscheme=ldaps`. To use encrypted LDAP connections using the `StartTLS` operation, use the normal URL scheme `ldap` and specify the `ldaptls` option in addition to `ldapurl`.

For non-anonymous binds, `ldapbinddn` and `ldapbindpasswd` must be specified as separate options.

LDAP URLs are currently only supported with OpenLDAP, not on Windows.

It is an error to mix configuration options for simple bind with options for search+bind.

When using search+bind mode, the search can be performed using a single attribute specified with `ldapsearchattribute`, or using a custom search filter specified with `ldapsearchfilter`. Specifying `ldapsearchattribute=foo` is equivalent to specifying `ldapsearchfilter="(foo=$username)"`. If neither option is specified the default is `ldapsearchattribute=uid`.

If Postgres Pro was compiled with OpenLDAP as the LDAP client library, the `ldapserver` setting may be omitted. In that case, a list of host names and ports is looked up via [RFC 2782](#) DNS SRV records. The name `_ldap._tcp.DOMAIN` is looked up, where `DOMAIN` is extracted from `ldapbasedn`.

Here is an example for a simple-bind LDAP configuration:

```
host ... ldap ldapserver=ldap.example.net ldapprefix="cn=" ldapsuffix=", dc=example, dc=net"
```

When a connection to the database server as database user `someuser` is requested, Postgres Pro will attempt to bind to the LDAP server using the DN `cn=someuser, dc=example, dc=net` and the password provided by the client. If that connection succeeds, the database access is granted.

Here is an example for a search+bind configuration:

```
host ... ldap ldapserver=ldap.example.net ldapbasedn="dc=example, dc=net"
      ldapsearchattribute=uid
```

When a connection to the database server as database user `someuser` is requested, Postgres Pro will attempt to bind anonymously (since `ldapbinddn` was not specified) to the LDAP server, perform a search for `(uid=someuser)` under the specified base DN. If an entry is found, it will then attempt to bind using that found information and the password supplied by the client. If that second bind succeeds, the database access is granted.

Here is the same search+bind configuration written as a URL:

```
host ... ldap ldapurl="ldap://ldap.example.net/dc=example,dc=net?uid?sub"
```

Some other software that supports authentication against LDAP uses the same URL format, so it will be easier to share the configuration.

Here is an example for a search+bind configuration that uses `ldapsearchfilter` instead of `ldapsearchattribute` to allow authentication by user ID or email address:

```
host ... ldap ldapserver=ldap.example.net ldapbasedn="dc=example, dc=net"
      ldapsearchfilter="( | (uid=$username) (mail=$username) )"
```

Here is an example for a search+bind configuration that uses DNS SRV discovery to find the host name(s) and port(s) for the LDAP service for the domain name `example.net`:

```
host ... ldap ldapbasedn="dc=example,dc=net"
```

Tip

Since LDAP often uses commas and spaces to separate the different parts of a DN, it is often necessary to use double-quoted parameter values when configuring LDAP options, as shown in the examples.

20.10.1. Recommendations on Configuring Common LDAP Servers

The following requirements must be satisfied to configure LDAP authentication on a Postgres Pro server:

- An LDAP server must be accessible from the Postgres Pro server.
- The LDAP authentication method must be the first record in the `pg_hba.conf` configuration file on the Postgres Pro server.

- If LDAP authentication is configured with the `ldapbinddn` option, then the corresponding DN must already exist on the LDAP server.
- If LDAP authentication is configured with the `ldapsrchfilter` option to filter users by checking their group membership, then the corresponding group must already exist on the LDAP server.

While some LDAP servers (like FreeIPA) allow anonymous binding in the search+bind mode by default, others (like ALD Pro) do not. It is not recommended to allow anonymous binding on any LDAP server as it may result in security issues.

Note

Postgres Pro can authenticate users via ALD Pro 2.4.0 or higher.

Examples of configuration records listed below are compatible with FreeIPA and ALD Pro LDAP servers.

Here is an example of a simple-bind LDAP configuration:

```
host ... ldap ldapserver=ldap.example.net \  
          ldapport=389 \  
          ldapbasedn="cn=users,cn=accounts,dc=example,dc=net"
```

Here is an example of a search+bind LDAP configuration:

```
host ... ldap ldapserver=ldap.example.net \  
          ldapport=389 \  
          ldapbasedn="cn=users,cn=accounts,dc=example,dc=net" \  
          ldapbinddn="uid=pgpro,cn=users,cn=accounts,dc=example,dc=net" \  
          ldapbindpasswd="<pgpro password>" \  
          ldapsrchattribute=uid
```

When a database connection is requested as user `someuser`, Postgres Pro will attempt to bind to the LDAP server as `pgpro` user, search for `(uid=someuser)` under the specified base DN. If a matching record is found, the database server will then attempt to bind using that found information and the password supplied by the client. If that second connection succeeds, the database access is granted.

Here is an example for a search+bind LDAP configuration that uses `ldapsrchfilter` to filter users who are members of the `pgpro_access` group:

```
host ... ldap ldapserver=ldap.example.net \  
          ldapport=389 \  
          ldapbasedn="cn=users,cn=accounts,dc=example,dc=net" \  
          ldapbinddn="uid=pgpro,cn=users,cn=accounts,dc=example,dc=net" \  
          ldapbindpasswd="<pgpro password>" \  
          ldapsrchfilter="(&(objectClass=person)(uid=$username)  
(memberOf=cn=pgpro_access,cn=groups,cn=accounts,dc=example,dc=net))"
```

When a database connection is requested as user `someuser`, Postgres Pro will attempt to bind to the LDAP server as `pgpro` user, search for `(uid=someuser)` under the specified base DN according to the specified search filter. In this example, a user must be a member of the `pgpro_access` group according to the filter. If a matching record is found, the database server will then attempt to bind using that found information and the password supplied by the client. If that second connection succeeds, the database access is granted.

20.11. RADIUS Authentication

This authentication method operates similarly to `password` except that it uses RADIUS as the password verification method. RADIUS is used only to validate the user name/password pairs. Therefore the user must already exist in the database before RADIUS can be used for authentication.

When using RADIUS authentication, an Access Request message will be sent to the configured RADIUS server. This request will be of type `Authenticate Only`, and include parameters for `user name`, `password` (encrypted) and `NAS Identifier`. The request will be encrypted using a secret shared with the server. The RADIUS server will respond to this request with either `Access Accept` or `Access Reject`. There is no support for RADIUS accounting.

Multiple RADIUS servers can be specified, in which case they will be tried sequentially. If a negative response is received from a server, the authentication will fail. If no response is received, the next server in the list will be tried. To specify multiple servers, separate the server names with commas and surround the list with double quotes. If multiple servers are specified, the other RADIUS options can also be given as comma-separated lists, to provide individual values for each server. They can also be specified as a single value, in which case that value will apply to all servers.

The following configuration options are supported for RADIUS:

`radiusservers`

The DNS names or IP addresses of the RADIUS servers to connect to. This parameter is required.

`radiussecrets`

The shared secrets used when talking securely to the RADIUS servers. This must have exactly the same value on the Postgres Pro and RADIUS servers. It is recommended that this be a string of at least 16 characters. This parameter is required.

Note

The encryption vector used will only be cryptographically strong if Postgres Pro is built with support for OpenSSL. In other cases, the transmission to the RADIUS server should only be considered obfuscated, not secured, and external security measures should be applied if necessary.

`radiusports`

The port numbers to connect to on the RADIUS servers. If no port is specified, the default RADIUS port (1812) will be used.

`radiusidentifiers`

The strings to be used as `NAS Identifier` in the RADIUS requests. This parameter can be used, for example, to identify which database cluster the user is attempting to connect to, which can be useful for policy matching on the RADIUS server. If no identifier is specified, the default `postgresql` will be used.

If it is necessary to have a comma or whitespace in a RADIUS parameter value, that can be done by putting double quotes around the value, but it is tedious because two layers of double-quoting are now required. An example of putting whitespace into RADIUS secret strings is:

```
host ... radius radiusservers="server1,server2" radiussecrets="\"secret one\", \"secret two\""
```

20.12. Certificate Authentication

This authentication method uses SSL client certificates to perform authentication. It is therefore only available for SSL connections; see [Section 18.9.2](#) for SSL configuration instructions. When using this authentication method, the server will require that the client provide a valid, trusted certificate. No

password prompt will be sent to the client. The `cn` (Common Name) attribute of the certificate will be compared to the requested database user name, and if they match the login will be allowed. User name mapping can be used to allow `cn` to be different from the database user name.

The following configuration options are supported for SSL certificate authentication:

`map`

Allows for mapping between system and database user names. See [Section 20.2](#) for details.

It is redundant to use the `clientcert` option with `cert` authentication because `cert` authentication is effectively `trust` authentication with `clientcert=verify-full`.

20.13. PAM Authentication

This authentication method operates similarly to `password` except that it uses PAM (Pluggable Authentication Modules) as the authentication mechanism. The default PAM service name is `postgresql`. PAM is used only to validate user name/password pairs and optionally the connected remote host name or IP address. Therefore the user must already exist in the database before PAM can be used for authentication. For more information about PAM, please read the [Linux-PAM Page](#).

The following configuration options are supported for PAM:

`pamservice`

PAM service name.

`pam_use_hostname`

Determines whether the remote IP address or the host name is provided to PAM modules through the `PAM_RHOST` item. By default, the IP address is used. Set this option to 1 to use the resolved host name instead. Host name resolution can lead to login delays. (Most PAM configurations don't use this information, so it is only necessary to consider this setting if a PAM configuration was specifically created to make use of it.)

Note

If PAM is set up to read `/etc/shadow`, authentication will fail because the Postgres Pro server is started by a non-root user. However, this is not an issue when PAM is configured to use LDAP or other authentication methods.

20.14. BSD Authentication

This authentication method operates similarly to `password` except that it uses BSD Authentication to verify the password. BSD Authentication is used only to validate user name/password pairs. Therefore the user's role must already exist in the database before BSD Authentication can be used for authentication. The BSD Authentication framework is currently only available on OpenBSD.

BSD Authentication in Postgres Pro uses the `auth-postgresql` login type and authenticates with the `postgresql` login class if that's defined in `login.conf`. By default that login class does not exist, and Postgres Pro will use the default login class.

Note

To use BSD Authentication, the Postgres Pro user account (that is, the operating system user running the server) must first be added to the `auth` group. The `auth` group exists by default on OpenBSD systems.

20.15. Authentication Problems

Authentication failures and related problems generally manifest themselves through error messages like the following:

```
FATAL: no pg_hba.conf entry for host "123.123.123.123", user "andym", database
"testdb"
```

This is what you are most likely to get if you succeed in contacting the server, but it does not want to talk to you. As the message suggests, the server refused the connection request because it found no matching entry in its `pg_hba.conf` configuration file.

```
FATAL: password authentication failed for user "andym"
```

Messages like this indicate that you contacted the server, and it is willing to talk to you, but not until you pass the authorization method specified in the `pg_hba.conf` file. Check the password you are providing, or check your Kerberos or ident software if the complaint mentions one of those authentication types.

```
FATAL: user "andym" does not exist
```

The indicated database user name was not found.

```
FATAL: database "testdb" does not exist
```

The database you are trying to connect to does not exist. Note that if you do not specify a database name, it defaults to the database user name.

Tip

The server log might contain more information about an authentication failure than is reported to the client. If you are confused about the reason for a failure, check the server log.

Chapter 21. Database Roles

Postgres Pro manages database access permissions using the concept of *roles*. A role can be thought of as either a database user, or a group of database users, depending on how the role is set up. Roles can own database objects (for example, tables and functions) and can assign privileges on those objects to other roles to control who has access to which objects. Furthermore, it is possible to grant *membership* in a role to another role, thus allowing the member role to use privileges assigned to another role.

The concept of roles subsumes the concepts of “users” and “groups”. In PostgreSQL versions before 8.1, users and groups were distinct kinds of entities, but now there are only roles. Any role can act as a user, a group, or both.

This chapter describes how to create and manage roles. More information about the effects of role privileges on various database objects can be found in [Section 5.8](#).

21.1. Database Roles

Database roles are conceptually completely separate from operating system users. In practice it might be convenient to maintain a correspondence, but this is not required. Database roles are global across a database cluster installation (and not per individual database). To create a role use the [CREATE ROLE](#) SQL command:

```
CREATE ROLE name;
```

name follows the rules for SQL identifiers: either unadorned without special characters, or double-quoted. (In practice, you will usually want to add additional options, such as [LOGIN](#), to the command. More details appear below.) To remove an existing role, use the analogous [DROP ROLE](#) command:

```
DROP ROLE name;
```

For convenience, the programs [createuser](#) and [dropuser](#) are provided as wrappers around these SQL commands that can be called from the shell command line:

```
createuser name  
dropuser name
```

To determine the set of existing roles, examine the `pg_roles` system catalog, for example:

```
SELECT rolname FROM pg_roles;
```

or to see just those capable of logging in:

```
SELECT rolname FROM pg_roles WHERE rolcanlogin;
```

The [psql](#) program's `\du` meta-command is also useful for listing the existing roles.

In order to bootstrap the database system, a freshly initialized system always contains one predefined login-capable role. This role is always a “superuser”, and it will have the same name as the operating system user that initialized the database cluster with `initdb` unless a different name is specified. This role is often named `postgres`. In order to create more roles you first have to connect as this initial role.

Every connection to the database server is made using the name of some particular role, and this role determines the initial access privileges for commands issued in that connection. The role name to use for a particular database connection is indicated by the client that is initiating the connection request in an application-specific fashion. For example, the `psql` program uses the `-U` command line option to indicate the role to connect as. Many applications assume the name of the current operating system user by default (including `createuser` and `psql`). Therefore it is often convenient to maintain a naming correspondence between roles and operating system users.

The set of database roles a given client connection can connect as is determined by the client authentication setup, as explained in [Chapter 20](#). (Thus, a client is not limited to connect as the role matching

its operating system user, just as a person's login name need not match his or her real name.) Since the role identity determines the set of privileges available to a connected client, it is important to carefully configure privileges when setting up a multiuser environment.

21.2. Role Attributes

A database role can have a number of attributes that define its privileges and interact with the client authentication system.

login privilege

Only roles that have the `LOGIN` attribute can be used as the initial role name for a database connection. A role with the `LOGIN` attribute can be considered the same as a “database user”. To create a role with login privilege, use either:

```
CREATE ROLE name LOGIN;  
CREATE USER name;
```

(`CREATE USER` is equivalent to `CREATE ROLE` except that `CREATE USER` includes `LOGIN` by default, while `CREATE ROLE` does not.)

superuser status

A database superuser bypasses all permission checks, except the right to log in. This is a dangerous privilege and should not be used carelessly; it is best to do most of your work as a role that is not a superuser. To create a new database superuser, use `CREATE ROLE name SUPERUSER`. You must do this as a role that is already a superuser.

database creation

A role must be explicitly given permission to create databases (except for superusers, since those bypass all permission checks). To create such a role, use `CREATE ROLE name CREATEDB`.

role creation

A role must be explicitly given permission to create more roles (except for superusers, since those bypass all permission checks). To create such a role, use `CREATE ROLE name CREATEROLE`. A role with `CREATEROLE` privilege can alter and drop roles which have been granted to the `CREATEROLE` user with the `ADMIN` option. Such a grant occurs automatically when a `CREATEROLE` user that is not a superuser creates a new role, so that by default, a `CREATEROLE` user can alter and drop the roles which they have created. Altering a role includes most changes that can be made using `ALTER ROLE`, including, for example, changing passwords. It also includes modifications to a role that can be made using the `COMMENT` and `SECURITY LABEL` commands.

However, `CREATEROLE` does not convey the ability to create `SUPERUSER` roles, nor does it convey any power over `SUPERUSER` roles that already exist. Furthermore, `CREATEROLE` does not convey the power to create `REPLICATION` users, nor the ability to grant or revoke the `REPLICATION` privilege, nor the ability to modify the role properties of such users. However, it does allow `ALTER ROLE ... SET` and `ALTER ROLE ... RENAME` to be used on `REPLICATION` roles, as well as the use of `COMMENT ON ROLE`, `SECURITY LABEL ON ROLE`, and `DROP ROLE`. Finally, `CREATEROLE` does not confer the ability to grant or revoke the `BYPASSRLS` privilege.

initiating replication

A role must explicitly be given permission to initiate streaming replication (except for superusers, since those bypass all permission checks). A role used for streaming replication must have `LOGIN` permission as well. To create such a role, use `CREATE ROLE name REPLICATION LOGIN`.

password

A password is only significant if the client authentication method requires the user to supply a password when connecting to the database. The `password` and `md5` authentication methods make use of

passwords. Database passwords are separate from operating system passwords. Specify a password upon role creation with `CREATE ROLE name PASSWORD 'string'`.

inheritance of privileges

A role inherits the privileges of roles it is a member of, by default. However, to create a role which does not inherit privileges by default, use `CREATE ROLE name NOINHERIT`. Alternatively, inheritance can be overridden for individual grants by using `WITH INHERIT TRUE` or `WITH INHERIT FALSE`.

bypassing row-level security

A role must be explicitly given permission to bypass every row-level security (RLS) policy (except for superusers, since those bypass all permission checks). To create such a role, use `CREATE ROLE name BYPASSRLS` as a superuser.

connection limit

Connection limit can specify how many concurrent connections a role can make. -1 (the default) means no limit. Specify connection limit upon role creation with `CREATE ROLE name CONNECTION LIMIT 'integer'`.

A role's attributes can be modified after creation with `ALTER ROLE`. See the reference pages for the [CREATE ROLE](#) and [ALTER ROLE](#) commands for details.

A role can also have role-specific defaults for many of the run-time configuration settings described in [Chapter 19](#). For example, if for some reason you want to disable index scans (hint: not a good idea) anytime you connect, you can use:

```
ALTER ROLE myname SET enable_indexscan TO off;
```

This will save the setting (but not set it immediately). In subsequent connections by this role it will appear as though `SET enable_indexscan TO off` had been executed just before the session started. You can still alter this setting during the session; it will only be the default. To remove a role-specific default setting, use `ALTER ROLE rolename RESET varname`. Note that role-specific defaults attached to roles without `LOGIN` privilege are fairly useless, since they will never be invoked.

When a non-superuser creates a role using the `CREATEROLE` privilege, the created role is automatically granted back to the creating user, just as if the bootstrap superuser had executed the command `GRANT created_user TO creating_user WITH ADMIN TRUE, SET FALSE, INHERIT FALSE`. Since a `CREATEROLE` user can only exercise special privileges with regard to an existing role if they have `ADMIN OPTION` on it, this grant is just sufficient to allow a `CREATEROLE` user to administer the roles they created. However, because it is created with `INHERIT FALSE, SET FALSE`, the `CREATEROLE` user doesn't inherit the privileges of the created role, nor can it access the privileges of that role using `SET ROLE`. However, since any user who has `ADMIN OPTION` on a role can grant membership in that role to any other user, the `CREATEROLE` user can gain access to the created role by simply granting that role back to themselves with the `INHERIT` and/or `SET` options. Thus, the fact that privileges are not inherited by default nor is `SET ROLE` granted by default is a safeguard against accidents, not a security feature. Also note that, because this automatic grant is granted by the bootstrap superuser, it cannot be removed or changed by the `CREATEROLE` user; however, any superuser could revoke it, modify it, and/or issue additional such grants to other `CREATEROLE` users. Whichever `CREATEROLE` users have `ADMIN OPTION` on a role at any given time can administer it.

21.3. Role Membership

It is frequently convenient to group users together to ease management of privileges: that way, privileges can be granted to, or revoked from, a group as a whole. In Postgres Pro this is done by creating a role that represents the group, and then granting *membership* in the group role to individual user roles.

To set up a group role, first create the role:

```
CREATE ROLE name;
```


Typically a role being used as a group would not have the `LOGIN` attribute, though you can set it if you wish.

Once the group role exists, you can add and remove members using the `GRANT` and `REVOKE` commands:

```
GRANT group_role TO role1, ... ;
REVOKE group_role FROM role1, ... ;
```

You can grant membership to other group roles, too (since there isn't really any distinction between group roles and non-group roles). The database will not let you set up circular membership loops. Also, it is not permitted to grant membership in a role to `PUBLIC`.

The members of a group role can use the privileges of the role in two ways. First, member roles that have been granted membership with the `SET` option can do `SET ROLE` to temporarily “become” the group role. In this state, the database session has access to the privileges of the group role rather than the original login role, and any database objects created are considered owned by the group role not the login role. Second, member roles that have been granted membership with the `INHERIT` option automatically have use of the privileges of those directly or indirectly a member of, though the chain stops at memberships lacking the inherit option. As an example, suppose we have done:

```
CREATE ROLE joe LOGIN;
CREATE ROLE admin;
CREATE ROLE wheel;
CREATE ROLE island;
GRANT admin TO joe WITH INHERIT TRUE;
GRANT wheel TO admin WITH INHERIT FALSE;
GRANT island TO joe WITH INHERIT TRUE, SET FALSE;
```

Immediately after connecting as role `joe`, a database session will have use of privileges granted directly to `joe` plus any privileges granted to `admin` and `island`, because `joe` “inherits” those privileges. However, privileges granted to `wheel` are not available, because even though `joe` is indirectly a member of `wheel`, the membership is via `admin` which was granted using `WITH INHERIT FALSE`. After:

```
SET ROLE admin;
```

the session would have use of only those privileges granted to `admin`, and not those granted to `joe` or `island`. After:

```
SET ROLE wheel;
```

the session would have use of only those privileges granted to `wheel`, and not those granted to either `joe` or `admin`. The original privilege state can be restored with any of:

```
SET ROLE joe;
SET ROLE NONE;
RESET ROLE;
```

Note

The `SET ROLE` command always allows selecting any role that the original login role is directly or indirectly a member of, provided that there is a chain of membership grants each of which has `SET TRUE` (which is the default). Thus, in the above example, it is not necessary to become `admin` before becoming `wheel`. On the other hand, it is not possible to become `island` at all; `joe` can only access those privileges via inheritance.

Note

In the SQL standard, there is a clear distinction between users and roles, and users do not automatically inherit privileges while roles do. This behavior can be obtained in Postgres Pro by giving roles being used as SQL roles the `INHERIT` attribute, while giving roles being used as SQL users the `NOINHERIT` attribute. However, Postgres Pro defaults to giving all roles the `INHERIT` attribute,

for backward compatibility with pre-8.1 releases in which users always had use of permissions granted to groups they were members of.

The role attributes `LOGIN`, `SUPERUSER`, `CREATEDB`, and `CREATEROLE` can be thought of as special privileges, but they are never inherited as ordinary privileges on database objects are. You must actually `SET ROLE` to a specific role having one of these attributes in order to make use of the attribute. Continuing the above example, we might choose to grant `CREATEDB` and `CREATEROLE` to the `admin` role. Then a session connecting as role `joe` would not have these privileges immediately, only after doing `SET ROLE admin`.

To destroy a group role, use `DROP ROLE`:

```
DROP ROLE name;
```

Any memberships in the group role are automatically revoked (but the member roles are not otherwise affected).

21.4. Dropping Roles

Because roles can own database objects and can hold privileges to access other objects, dropping a role is often not just a matter of a quick `DROP ROLE`. Any objects owned by the role must first be dropped or reassigned to other owners; and any permissions granted to the role must be revoked.

Ownership of objects can be transferred one at a time using `ALTER` commands, for example:

```
ALTER TABLE bobs_table OWNER TO alice;
```

Alternatively, the `REASSIGN OWNED` command can be used to reassign ownership of all objects owned by the role-to-be-dropped to a single other role. Because `REASSIGN OWNED` cannot access objects in other databases, it is necessary to run it in each database that contains objects owned by the role. (Note that the first such `REASSIGN OWNED` will change the ownership of any shared-across-databases objects, that is databases or tablespaces, that are owned by the role-to-be-dropped.)

Once any valuable objects have been transferred to new owners, any remaining objects owned by the role-to-be-dropped can be dropped with the `DROP OWNED` command. Again, this command cannot access objects in other databases, so it is necessary to run it in each database that contains objects owned by the role. Also, `DROP OWNED` will not drop entire databases or tablespaces, so it is necessary to do that manually if the role owns any databases or tablespaces that have not been transferred to new owners.

`DROP OWNED` also takes care of removing any privileges granted to the target role for objects that do not belong to it. Because `REASSIGN OWNED` does not touch such objects, it's typically necessary to run both `REASSIGN OWNED` and `DROP OWNED` (in that order!) to fully remove the dependencies of a role to be dropped.

In short then, the most general recipe for removing a role that has been used to own objects is:

```
REASSIGN OWNED BY doomed_role TO successor_role;
DROP OWNED BY doomed_role;
-- repeat the above commands in each database of the cluster
DROP ROLE doomed_role;
```

When not all owned objects are to be transferred to the same successor owner, it's best to handle the exceptions manually and then perform the above steps to mop up.

If `DROP ROLE` is attempted while dependent objects still remain, it will issue messages identifying which objects need to be reassigned or dropped.

21.5. Predefined Roles

Postgres Pro provides a set of predefined roles that provide access to certain, commonly needed, privileged capabilities and information. Administrators (including roles that have the `CREATEROLE` privilege) can `GRANT` these roles to users and/or other roles in their environment, providing those users with access to the specified capabilities and information.

The predefined roles are described in [Table 21.1](#). Note that the specific permissions for each of the roles may change in the future as additional capabilities are added. Administrators should monitor the release notes for changes.

Table 21.1. Predefined Roles

Role	Allowed Access
pg_read_all_data	Read all data (tables, views, sequences), as if having <code>SELECT</code> rights on those objects, and <code>USAGE</code> rights on all schemas, even without having it explicitly. This role does not have the role attribute <code>BYPASSRLS</code> set. If RLS is being used, an administrator may wish to set <code>BYPASSRLS</code> on roles which this role is GRANTED to.
pg_write_all_data	Write all data (tables, views, sequences), as if having <code>INSERT</code> , <code>UPDATE</code> , and <code>DELETE</code> rights on those objects, and <code>USAGE</code> rights on all schemas, even without having it explicitly. This role does not have the role attribute <code>BYPASSRLS</code> set. If RLS is being used, an administrator may wish to set <code>BYPASSRLS</code> on roles which this role is GRANTED to.
pg_read_all_settings	Read all configuration variables, even those normally visible only to superusers.
pg_read_all_stats	Read all <code>pg_stat_*</code> views and use various statistics related extensions, even those normally visible only to superusers.
pg_stat_scan_tables	Execute monitoring functions that may take <code>ACCESS SHARE</code> locks on tables, potentially for a long time.
pg_monitor	Read/execute various monitoring views and functions. This role is a member of <code>pg_read_all_settings</code> , <code>pg_read_all_stats</code> and <code>pg_stat_scan_tables</code> .
pg_database_owner	None. Membership consists, implicitly, of the current database owner.
pg_signal_backend	Signal another backend to cancel a query or terminate its session.
pg_read_server_files	Allow reading files from any location the database can access on the server with <code>COPY</code> and other file-access functions.
pg_write_server_files	Allow writing to files in any location the database can access on the server with <code>COPY</code> and other file-access functions.
pg_execute_server_program	Allow executing programs on the database server as the user the database runs as with <code>COPY</code> and other functions which allow executing a server-side program.
pg_checkpoint	Allow executing the <code>CHECKPOINT</code> command.
pg_maintain	Allow executing <code>VACUUM</code> , <code>ANALYZE</code> , <code>CLUSTER</code> , <code>REFRESH MATERIALIZED VIEW</code> , <code>REINDEX</code> , and <code>LOCK TABLE</code> on all relations, as if having <code>MAINTAIN</code> rights on those objects, even without having it explicitly.
pg_use_reserved_connections	Allow use of connection slots reserved via <code>reserved_connections</code> .
pg_create_subscription	Allow users with <code>CREATE</code> permission on the database to issue <code>CREATE SUBSCRIPTION</code> .
pg_create_tablespace	Execute the <code>CREATE TABLESPACE</code> command without superuser rights.
pg_manage_profiles	Execute the <code>CREATE PROFILE</code> , <code>ALTER PROFILE</code> , and <code>DROP PROFILE</code> commands without superuser rights.

The `pg_monitor`, `pg_read_all_settings`, `pg_read_all_stats` and `pg_stat_scan_tables` roles are intended to allow administrators to easily configure a role for the purpose of monitoring the database

server. They grant a set of common privileges allowing the role to read various useful configuration settings, statistics and other system information normally restricted to superusers.

The `pg_database_owner` role has one implicit, situation-dependent member, namely the owner of the current database. Like any role, it can own objects or receive grants of access privileges. Consequently, once `pg_database_owner` has rights within a template database, each owner of a database instantiated from that template will exercise those rights. `pg_database_owner` cannot be a member of any role, and it cannot have non-implicit members. Initially, this role owns the `public` schema, so each database owner governs local use of the schema.

The `pg_signal_backend` role is intended to allow administrators to enable trusted, but non-superuser, roles to send signals to other backends. Currently this role enables sending of signals for canceling a query on another backend or terminating its session. A user granted this role cannot however send signals to a backend owned by a superuser. See [Section 9.28.2](#).

The `pg_read_server_files`, `pg_write_server_files` and `pg_execute_server_program` roles are intended to allow administrators to have trusted, but non-superuser, roles which are able to access files and run programs on the database server as the user the database runs as. As these roles are able to access any file on the server file system, they bypass all database-level permission checks when accessing files directly and they could be used to gain superuser-level access, therefore great care should be taken when granting these roles to users.

Care should be taken when granting these roles to ensure they are only used where needed and with the understanding that these roles grant access to privileged information.

Administrators can grant access to these roles to users using the `GRANT` command, for example:

```
GRANT pg_signal_backend TO admin_user;
```

21.6. Function Security

Functions, triggers and row-level security policies allow users to insert code into the backend server that other users might execute unintentionally. Hence, these mechanisms permit users to “Trojan horse” others with relative ease. The strongest protection is tight control over who can define objects. Where that is infeasible, write queries referring only to objects having trusted owners. Remove from `search_path` any schemas that permit untrusted users to create objects.

Functions run inside the backend server process with the operating system permissions of the database server daemon. If the programming language used for the function allows unchecked memory accesses, it is possible to change the server's internal data structures. Hence, among many other things, such functions can circumvent any system access controls. Function languages that allow such access are considered “untrusted”, and Postgres Pro allows only superusers to create functions written in those languages.

21.7. Roles in Distributed System

21.7.1. Managing Users and Roles

Users and roles in a Postgres Pro Shardman cluster are usual PostgreSQL users and roles. You can manage them separately on each server or globally, using broadcast DDL. Postgres Pro Shardman also uses concepts of *global users* and *global roles*. And only the *global users* (or roles) can create and own other Postgres Pro Shardman cluster-wide objects, such as sharded or global tables. Operations on such users are always performed on all replication groups simultaneously. For example, when you include a global role in some other role or drop it, this operation will be performed on all replication groups.

You can create a global user with a `CREATE USER ... IN ROLE global` statement, for example:

```
CREATE USER someuser ENCRYPTED PASSWORD 'somepass' IN ROLE global;
```

When a global user is created, Postgres Pro Shardman automatically creates user mappings on all replication groups and grants this user with access to all foreign servers corresponding to existing replica-

tion groups. Therefore, when you create a global user, you need to specify either a cleartext password, so that it can be saved in a user mapping, or no password at all. A passwordless global user or role is unable to access foreign servers, but you can use such a role to accumulate some permissions and grant it to different users. You can also set a password for a passwordless global user later.

Global users can be created only by user with `CREATEROLE` permission on all cluster nodes.

`ALTER` and `DROP` statements for global users are broadcasted to all replication groups. When a role is granted to a global user, this operation is also broadcasted. Renaming a global user is not supported since this invalidates md5/scram-sha-256 passwords stored in user mappings.

The list of global users is stored in the `shardman.users` table.

The role specified in `PgSuUsername` (usually, `postgres`) is also created as global user during cluster initialization. However, the role specified in `PgReplUsername` is created as local user on each replication group.

The role `global` is reserved and cannot be used directly in a Shardman cluster. Note that 'global' is not a really defined role but just a reserved word.

21.7.2. Managing Permissions for Sharded Tables

In Postgres Pro Shardman, a sharded table is basically a partitioned table where partitions are either local shards or foreign tables referencing shards in other replication groups.

Permissions granted on a sharded table are broadcasted to all replication groups and to all partitions of the table.

When a new replication group is added to a cluster, `shardmanctl` copies the schema from a random existing replication group to the new one. It also creates a foreign server for the new replication group on all existing replication groups and recreates foreign servers on new replication groups. Permissions for the created foreign servers and user mappings are copied from a random foreign server in an existing replication group. In the new replication group, for each partition of the sharded table `shardmanctl` creates a foreign table referencing the existing shard and replaces the partition with this foreign table. Later some of these foreign tables can be replaced by real tables. This happens during the `shardmanctl nodes add rebalance` stage when rebalance is enabled. Data for these partitions is transferred from existing nodes using logical replication. When `shardmanctl` creates tables (or foreign tables), it copies permissions from the parent table. The parent table must already have correct permissions since they were copied from an existing replication group.

Chapter 22. Managing Databases

Every instance of a running Postgres Pro server manages one or more databases. Databases are therefore the topmost hierarchical level for organizing SQL objects (“database objects”). This chapter describes the properties of databases, and how to create, manage, and destroy them.

The managing principles described in this chapter are applied for regular and distributed systems, yet only for local sharded DBs as, for now, sharding works only for a database named `postgres` (default), creating other databases is in development.

22.1. Postgres Pro Shardman Overview

Postgres Pro Shardman is a distributed system, which has both advantages and disadvantages inherent to such systems. Besides, Postgres Pro Shardman is primarily designed for OLTP load. OLAP queries to Postgres Pro Shardman are also possible, but only pretty simple of them. If you want to load an OLTP system with OLAP functionality, bear in mind that the lists of analytic and aggregate SQL functions to be sent to other shards for execution are highly limited.

Postgres Pro Shardman does not have a single DB server, instead there is a pool of shard servers with equal privileges.

Special attention should be paid to type casts in queries because inclusion of a type casting function in a query condition can make it impossible to be pushed down to a remote server.

For a detailed description of the main terms, refer to [Glossary](#).

When connecting to the database server, a client must specify the database name in its connection request. It is not possible to access more than one database per connection. However, clients can open multiple connections to the same database. Database-level security has two components: access control (see [Section 20.1](#)), managed at the connection level, and authorization control (see [Section 5.8](#)), managed via the grant system.

As Postgres Pro Shardman can store several databases, each one in its shard, it is advised to consider carefully whether the benefits outweigh the risks and limitations. In particular, the impact that having a shared WAL (see [Chapter 28](#)) has on backup and recovery options. While individual databases in the cluster are isolated when considered from the user's perspective, they are closely bound from the database administrator's point-of-view.

Databases are created with the `CREATE DATABASE` command (see [Section 22.3](#)) and destroyed with the `DROP DATABASE` command (see [Section 22.6](#)). To determine the set of existing databases, examine the `pg_database` system catalog. for example

```
SELECT datname FROM pg_database;
```

The `psql` program's `\l` meta-command and `-l` command-line option are also useful for listing the existing databases.

Note

The SQL standard calls databases “catalogs”, but there is no difference in practice.

Taking into account the above features and limitations of the RDBMS, we will provide two simple examples of the transition from a regular to a distributed database schema.

22.2. Connecting and Working with a Postgres Pro Shardman Cluster

The cluster considered consists of four shards. This is how the data partitions of the main sharded table are distributed across shards.

For the `ticket_no` sharding key:

- `tickets_0` — shard-1 (cluster node `node1`)
- `tickets_1` — shard-2 (cluster node `node2`)
- `tickets_2` — shard-3 (cluster node `node3`)
- `tickets_3` — shard-4 (cluster node `node4`)

For the `book_ref` sharding key:

- `bookings_0` — shard-1 (cluster node `node1`)
- `bookings_1` — shard-2 (cluster node `node2`)
- `bookings_2` — shard-3 (cluster node `node3`)
- `bookings_3` — shard-4 (cluster node `node4`)

The examples below are provided for the `book_ref` sharding key, but the code in the subsections is suitable for the `ticket_no` sharding key.

Do not treat this code as optimal or use it in a production environment. It only shows how to implement creation of a connection pool to work with a Postgres Pro Shardman cluster.

What is common for all the examples is the cluster connection string, which must contain node names, TCP port numbers, user name and password, database name for connection and a set of session parameters.

You can get this string using the [shardmanctl](#) utility. In the simplest case, the string looks like this:

```
$ shardmanctl getconnstr

dbname=postgres host=node1,node2,node3,node4 port=5432,5432,5432,5432
```

You can get this string to connect to cluster nodes or to create the connection pool in applications.

22.2.1. SQL

A few convenient functions and views are implemented in Postgres Pro Shardman that add cluster observability by:

- Listing global tables
- Listing sharded tables
- Listing global sequences
- Finding the shard number from the value of the sharding key
- Performing `ANALYZE` for all the global and sharded tables in the cluster

22.2.1.1. Listing Global Tables

To display all global tables in the cluster, use the `shardman.global_tables` view:

```
postgres=# select
    relname as table_name,
    nspname as schema
from shardman.global_tables;

table_name | schema
-----+-----
aircrafts  | bookings
```



```
seats      | bookings
airports   | bookings
flights    | bookings
(4 rows)
```

22.2.1.2. Listing Sharded Tables

To display information on all the sharded tables in the cluster, query the `shardman.sharded_tables` view as follows:

```
postgres=# select
    relname as table_name,
    nparts  as partitions,
    colocated_with::oid::regclass::text as colocated_with,
    nspname as schema
from shardman.sharded_tables;
```

table_name	partitions	colocated_with	schema
bookings	4		bookings
ticket_flights	4	bookings	bookings
tickets	4	bookings	bookings
boarding_passes	4	bookings	bookings

(4 rows)

22.2.1.3. Listing Global Sequences

To display all the global sequences in the cluster, use the `shardman.sequence` view:

```
postgres=# select
    seqns  as schema,
    seqname as sequence_name,
    seqmin  as min_value,
    seqmax  as max_value,
    seqblk  as bulk_size
from shardman.sequence;
```

schema	sequence_name	min_value	max_value	bulk_size
bookings	flights_flight_id_seq	262145	9223372036854775807	65536

(1 rows)

22.2.1.4. Finding the Shard Number from the Sharding Key Value

To display the name of the partition that contains data and the replication group name, call the `shardman.get_partition_for_value()` function. For example, for `book_ref = 0369E5`:

```
postgres=# select * from shardman.get_partition_for_value(
    'bookings'::regclass,
    '0369E5'::character(6));
```

rgid	local_nspname	local_relname	remote_nspname	remote_relname
1	bookings	bookings_0	bookings	bookings_0

This output shows that the data is in the `bookings_0` partition of the `bookings` table and is located on the node where the query was executed.

Let's create a query to display the name of the server where the partition with data is located. If we connect to the server that contains the partition, the server name is displayed as “current server”. If the data is on a different server, the hostname of the shard master is displayed:

```
SELECT p.rgid,
       local_relname AS partition_name,
       CASE
         WHEN r.srvid IS NULL THEN 'current server'
         ELSE (SELECT (SELECT split_part(kv, '=', 2)
                        FROM (SELECT unnest(fs.srvoptions) as kv) x
                        WHERE split_part(kv, '=', 1) = 'host')
              FROM shardman.repgroups rg
              JOIN pg_catalog.pg_foreign_server AS fs ON fs.oid = rg.srvid
              WHERE rg.id = p.rgid)
         END
         AS server_name
FROM shardman.get_partition_for_value('bookings'::regclass, '0369E5'::character(6))
p
      JOIN shardman.repgroups AS r ON
      r.id = p.rgid;
```

rgid	partition_name	server_name
1	bookings_0	current server

(1 row)

Execution of this query with another value of the sharding key, 0369E6, produces the output:

rgid	partition_name	server_name
4	bookings_3_fdw	node4

(1 row)

It is clear that the partition is on the node4 node.

Also note that the `shardman.rgid` parameter allows you to find the number of the node with the connection session. To do this, execute the query:

```
SELECT pg_catalog.current_setting('shardman.rgid');
```

You can use this value to determine the location of connection sessions for queries like discussed in this section.

The `shardman.get_partition_for_value()` is mainly designed for administration purposes, to better understand the data topology.

As a rule, do not use administration functions when writing SQL code for data access.

22.2.1.5. Understanding How Partitions of Sharded Tables Are Distributed Across Shards

You can get the list of all sharded tables in the `bookings` schema, together with the number of partitions and their distribution across servers (shards) from Postgres Pro Shardman metadata on any cluster node.

Consider the following query:

```
SELECT p.rel::regclass::text AS table_name,
       p.pnum,
       p.rgid,
```

```

        r.srvid,
        fs.srvname
FROM shardman.parts p
JOIN shardman.repgroups r
ON p.rgid = r.id
LEFT OUTER JOIN pg_foreign_server fs
    ON r.srvid = fs.oid;

```

To learn how the data is distributed, let's combine this query with a subquery from [Section 22.2.1.4](#):

```

SELECT p.rel::regclass AS table_name,
       st.nparts AS total_parts,
       p.pnum AS num_part,
       CASE
           WHEN r.srvid IS NULL THEN 'connected server'
       ELSE
           (SELECT split_part(kv, '=', 2)
            FROM (SELECT unnest(fs.srvoptions) AS kv) x
            WHERE split_part(kv, '=', 1) = 'host')
       END AS server_name
FROM shardman.parts p
JOIN shardman.repgroups r
    ON p.rgid = r.id
LEFT JOIN shardman.sharded_tables st
    ON p.rel = st.rel
LEFT JOIN pg_foreign_server fs
    ON r.srvid = fs.oid
WHERE st.nspname = 'bookings'
ORDER BY table_name, num_part, server_name;

```

The output format is the table name, number of table partitions, partition number and server name:

table_name	total_parts	num_part	server_name
bookings.bookings	4	0	connected server
bookings.bookings	4	1	node2
bookings.bookings	4	2	node3
bookings.bookings	4	3	node4
bookings.ticket_flights	4	0	connected server
bookings.ticket_flights	4	1	node2
bookings.ticket_flights	4	2	node3
bookings.ticket_flights	4	3	node4
bookings.tickets	4	0	connected server
bookings.tickets	4	1	node2
bookings.tickets	4	2	node3
bookings.tickets	4	3	node4
bookings.boarding_passes	4	0	connected server
bookings.boarding_passes	4	1	node2
bookings.boarding_passes	4	2	node3
bookings.boarding_passes	4	3	node4

22.2.1.6. Collecting Statistics

To collect statistics for sharded and global tables, call the `shardman.global_analyze()` function. This function first collects statistics for all local partitions of sharded tables on each node and then broadcasts this statistics to other nodes. For a global table, the function first collects statistics on a certain node and then the statistics is broadcast to all the other nodes.

22.3. Creating a Database

In this example, the table `pgbench_branches` is created, as well as colocated tables `pgbench_accounts` and `pgbench_history`. Each partition of the `pgbench_history` table is additionally subpartitioned by range.

```
CREATE TABLE pgbench_branches (  
    bid integer NOT NULL PRIMARY KEY,  
    bbalance integer,  
    filler character(88)  
)  
WITH (distributed_by = 'bid',  
      num_parts = 8);  
CREATE TABLE pgbench_accounts (  
    aid integer NOT NULL,  
    bid integer,  
    abalance integer,  
    filler character(84),  
    PRIMARY KEY (bid, aid)  
)  
WITH (distributed_by = 'bid',  
      num_parts = 8,  
      colocate_with = 'pgbench_branches');  
CREATE TABLE public.pgbench_history (  
    tid integer,  
    bid integer,  
    aid integer,  
    delta integer,  
    mtime timestamp without time zone,  
    filler character(22)  
)  
WITH (distributed_by = 'bid',  
      colocate_with = 'pgbench_branches',  
      partition_by = 'mtime',  
      partition_bounds =  
          $${{minvalue, '2021-01-01 00:00'}},{ '2021-01-01 00:00', '2022-01-01 00:00'}},  
      {'2022-01-01 00:00', maxvalue}}$$  
);
```

22.4. Template Databases

`CREATE DATABASE` actually works by copying an existing database. By default, it copies the standard system database named `template1`. Thus that database is the “template” from which new databases are made. If you add objects to `template1`, these objects will be copied into subsequently created user databases. This behavior allows site-local modifications to the standard set of objects in databases. For example, if you install the procedural language PL/Perl in `template1`, it will automatically be available in user databases without any extra action being taken when those databases are created.

However, `CREATE DATABASE` does not copy database-level `GRANT` permissions attached to the source database. The new database has default database-level permissions.

There is a second standard system database named `template0`. This database contains the same data as the initial contents of `template1`, that is, only the standard objects predefined by your version of Postgres Pro. `template0` should never be changed after the database cluster has been initialized. By instructing `CREATE DATABASE` to copy `template0` instead of `template1`, you can create a “pristine” user database (one where no user-defined objects exist and where the system objects have not been altered) that contains none of the site-local additions in `template1`. This is particularly handy when restoring a `pg_dump` dump: the dump script should be restored in a pristine database to ensure that one recreates

the correct contents of the dumped database, without conflicting with objects that might have been added to `template1` later on.

Another common reason for copying `template0` instead of `template1` is that new encoding and locale settings can be specified when copying `template0`, whereas a copy of `template1` must use the same settings it does. This is because `template1` might contain encoding-specific or locale-specific data, while `template0` is known not to.

To create a database by copying `template0`, use:

```
CREATE DATABASE dbname TEMPLATE template0;
```

from the SQL environment, or:

```
createdb -T template0 dbname
```

from the shell.

It is possible to create additional template databases, and indeed one can copy any database in a cluster by specifying its name as the template for `CREATE DATABASE`. It is important to understand, however, that this is not (yet) intended as a general-purpose “COPY DATABASE” facility. The principal limitation is that no other sessions can be connected to the source database while it is being copied. `CREATE DATABASE` will fail if any other connection exists when it starts; during the copy operation, new connections to the source database are prevented.

Two useful flags exist in `pg_database` for each database: the columns `datistemplate` and `dataallowconn`. `datistemplate` can be set to indicate that a database is intended as a template for `CREATE DATABASE`. If this flag is set, the database can be cloned by any user with `CREATEDB` privileges; if it is not set, only superusers and the owner of the database can clone it. If `dataallowconn` is false, then no new connections to that database will be allowed (but existing sessions are not terminated simply by setting the flag false). The `template0` database is normally marked `dataallowconn = false` to prevent its modification. Both `template0` and `template1` should always be marked with `datistemplate = true`.

Note

`template1` and `template0` do not have any special status beyond the fact that the name `template1` is the default source database name for `CREATE DATABASE`. For example, one could drop `template1` and recreate it from `template0` without any ill effects. This course of action might be advisable if one has carelessly added a bunch of junk in `template1`. (To delete `template1`, it must have `pg_database.datistemplate = false`.)

The `postgres` database is also created when a database cluster is initialized. This database is meant as a default database for users and applications to connect to. It is simply a copy of `template1` and can be dropped and recreated if necessary.

22.5. Database Configuration

Recall from [Chapter 19](#) that the Postgres Pro server provides a large number of run-time configuration variables. You can set database-specific default values for many of these settings.

For example, if for some reason you want to disable the GEQO optimizer for a given database, you'd ordinarily have to either disable it for all databases or make sure that every connecting client is careful to issue `SET geqo TO off`. To make this setting the default within a particular database, you can execute the command:

```
ALTER DATABASE mydb SET geqo TO off;
```

This will save the setting (but not set it immediately). In subsequent connections to this database it will appear as though `SET geqo TO off;` had been executed just before the session started. Note that users can still alter this setting during their sessions; it will only be the default. To undo any such setting, use `ALTER DATABASE dbname RESET varname.`

22.6. Destroying a Database

Databases are destroyed with the command [DROP DATABASE](#):

```
DROP DATABASE name;
```

Only the owner of the database, or a superuser, can drop a database. Dropping a database removes all objects that were contained within the database. The destruction of a database cannot be undone.

22.7. Tablespaces

Tablespaces in Postgres Pro allow database administrators to define locations in the file system where the files representing database objects can be stored. Once created, a tablespace can be referred to by name when creating database objects.

Creation of non-global tablespaces is not allowed by default. For more details, see [CREATE TABLESPACE](#).

By using tablespaces, an administrator can control the disk layout of a Postgres Pro installation. This is useful in at least two ways. First, if the partition or volume on which the cluster was initialized runs out of space and cannot be extended, a tablespace can be created on a different partition and used until the system can be reconfigured.

Second, tablespaces allow an administrator to use knowledge of the usage pattern of database objects to optimize performance. For example, an index which is very heavily used can be placed on a very fast, highly available disk, such as an expensive solid state device. At the same time a table storing archived data which is rarely used or not performance critical could be stored on a less expensive, slower disk system.

Warning

Even though located outside the main Postgres Pro data directory, tablespaces are an integral part of the database cluster and *cannot* be treated as an autonomous collection of data files. They are dependent on metadata contained in the main data directory, and therefore cannot be attached to a different database cluster or backed up individually. Similarly, if you lose a tablespace (file deletion, disk failure, etc.), the database cluster might become unreadable or unable to start. Placing a tablespace on a temporary file system like a RAM disk risks the reliability of the entire cluster.

To define a tablespace, use the [CREATE TABLESPACE](#) command, for example::

```
CREATE TABLESPACE fastspace LOCATION '/ssd1/postgresql/data';
```

The location must be an existing, empty directory that is owned by the Postgres Pro operating system user. All objects subsequently created within the tablespace will be stored in files underneath this directory. The location must not be on removable or transient storage, as the cluster might fail to function if the tablespace is missing or lost.

Note

There is usually not much point in making more than one tablespace per logical file system, since you cannot control the location of individual files within a logical file system. However, Postgres

Pro does not enforce any such limitation, and indeed it is not directly aware of the file system boundaries on your system. It just stores files in the directories you tell it to use.

Creation of the tablespace itself must be done as a database superuser or a user with the privileges of the `pg_create_tablespace` role, but after that you can allow ordinary database users to use it. To do that, grant them the `CREATE` privilege on it.

Tables, indexes, and entire databases can be assigned to particular tablespaces. To do so, a user with the `CREATE` privilege on a given tablespace must pass the tablespace name as a parameter to the relevant command. For example, the following creates a table in the tablespace `space1`:

```
CREATE TABLE foo(i int) TABLESPACE space1;
```

Alternatively, use the `default_tablespace` parameter:

```
SET default_tablespace = space1;  
CREATE TABLE foo(i int);
```

When `default_tablespace` is set to anything but an empty string, it supplies an implicit `TABLESPACE` clause for `CREATE TABLE` and `CREATE INDEX` commands that do not have an explicit one.

There is also a `temp_tablespaces` parameter, which determines the placement of temporary tables and indexes, as well as temporary files that are used for purposes such as sorting large data sets. This can be a list of tablespace names, rather than only one, so that the load associated with temporary objects can be spread over multiple tablespaces. A random member of the list is picked each time a temporary object is to be created.

The tablespace associated with a database is used to store the system catalogs of that database. Furthermore, it is the default tablespace used for tables, indexes, and temporary files created within the database, if no `TABLESPACE` clause is given and no other selection is specified by `default_tablespace` or `temp_tablespaces` (as appropriate). If a database is created without specifying a tablespace for it, it uses the same tablespace as the template database it is copied from.

Two tablespaces are automatically created when the database cluster is initialized. The `pg_global` tablespace is used only for shared system catalogs. The `pg_default` tablespace is the default tablespace of the `template1` and `template0` databases (and, therefore, will be the default tablespace for other databases as well, unless overridden by a `TABLESPACE` clause in `CREATE DATABASE`).

Once created, a tablespace can be used from any database, provided the requesting user has sufficient privilege. This means that a tablespace cannot be dropped until all objects in all databases using the tablespace have been removed.

To remove an empty tablespace, use the `DROP TABLESPACE` command.

To determine the set of existing tablespaces, examine the `pg_tablespace` system catalog, for example

```
SELECT spcname FROM pg_tablespace;
```

The `psql` program's `\db` meta-command is also useful for listing the existing tablespaces.

The directory `$PGDATA/pg_tblspc` contains symbolic links that point to each of the non-built-in tablespaces defined in the cluster. Although not recommended, it is possible to adjust the tablespace layout by hand by redefining these links. Under no circumstances perform this operation while the server is running. Note that in PostgreSQL 9.1 and earlier you will also need to update the `pg_tablespace` catalog with the new locations. (If you do not, `pg_dump` will continue to output the old tablespace locations.)

Chapter 23. Localization

This chapter describes the available localization features from the point of view of the administrator. Postgres Pro supports two localization facilities:

- Using the locale features of the operating system to provide locale-specific collation order, number formatting, translated messages, and other aspects. This is covered in [Section 23.1](#) and [Section 23.2](#).
- Providing a number of different character sets to support storing text in all kinds of languages, and providing character set translation between client and server. This is covered in [Section 23.3](#).

23.1. Locale Support

Locale support refers to an application respecting cultural preferences regarding alphabets, sorting, number formatting, etc. Postgres Pro uses the standard ISO C and POSIX locale facilities provided by the server operating system. For additional information refer to the documentation of your system.

23.1.1. Overview

Locale support is automatically initialized when a database cluster is created using `initdb`. `initdb` will initialize the database cluster with the locale setting of its execution environment by default, so if your system is already set to use the locale that you want in your database cluster then there is nothing else you need to do. If you want to use a different locale (or you are not sure which locale your system is set to), you can instruct `initdb` exactly which locale to use by specifying the `--locale` option. For example:

```
initdb --locale=ru_RU
```

This example for Unix systems sets the locale to Russian (`ru`) as spoken in Russia (`RU`). Other possibilities might include `en_US` (U.S. English) and `fr_CA` (French Canadian). If more than one character set can be used for a locale then the specifications can take the form *language_territory.codeset*. For example, `fr_BE.UTF-8` represents the French language (`fr`) as spoken in Belgium (`BE`), with a UTF-8 character set encoding.

What locales are available on your system under what names depends on what was provided by the operating system vendor and what was installed. On most Unix systems, the command `locale -a` will provide a list of available locales. Windows uses more verbose locale names, such as `German_Germany` or `Swedish_Sweden.1252`, but the principles are the same.

Occasionally it is useful to mix rules from several locales, e.g., use English collation rules but Spanish messages. To support that, a set of locale subcategories exist that control only certain aspects of the localization rules:

LC_COLLATE	String sort order
LC_CTYPE	Character classification (What is a letter? Its upper-case equivalent?)
LC_MESSAGES	Language of messages
LC_MONETARY	Formatting of currency amounts
LC_NUMERIC	Formatting of numbers
LC_TIME	Formatting of dates and times

The category names translate into names of `initdb` options to override the locale choice for a specific category. For instance, to set the locale to French Canadian, but use U.S. rules for formatting currency, use `initdb --locale=fr_CA --lc-monetary=en_US`.

If you want the system to behave as if it had no locale support, use the special locale name `C`, or equivalently `POSIX`.

Some locale categories must have their values fixed when the database is created. You can use different settings for different databases, but once a database is created, you cannot change them for that database anymore. `LC_COLLATE` and `LC_CTYPE` are these categories. They affect the sort order of indexes, so they must be kept fixed, or indexes on text columns would become corrupt. (But you can alleviate this restriction using collations, as discussed in [Section 23.2](#).) The default values for these categories are determined when `initdb` is run, and those values are used when new databases are created, unless specified otherwise in the `CREATE DATABASE` command.

The other locale categories can be changed whenever desired by setting the server configuration parameters that have the same name as the locale categories (see [Section 19.12.2](#) for details). The values that are chosen by `initdb` are actually only written into the configuration file `postgresql.conf` to serve as defaults when the server is started. If you remove these assignments from `postgresql.conf` then the server will inherit the settings from its execution environment.

Note that the locale behavior of the server is determined by the environment variables seen by the server, not by the environment of any client. Therefore, be careful to configure the correct locale settings before starting the server. A consequence of this is that if client and server are set up in different locales, messages might appear in different languages depending on where they originated.

Note

When we speak of inheriting the locale from the execution environment, this means the following on most operating systems: For a given locale category, say the collation, the following environment variables are consulted in this order until one is found to be set: `LC_ALL`, `LC_COLLATE` (or the variable corresponding to the respective category), `LANG`. If none of these environment variables are set then the locale defaults to `C`.

Some message localization libraries also look at the environment variable `LANGUAGE` which overrides all other locale settings for the purpose of setting the language of messages. If in doubt, please refer to the documentation of your operating system, in particular the documentation about `gettext`.

To enable messages to be translated to the user's preferred language, NLS must have been selected at build time (`configure --enable-nls`). All other locale support is built in automatically.

23.1.2. Behavior

The locale settings influence the following SQL features:

- Sort order in queries using `ORDER BY` or the standard comparison operators on textual data
- The `upper`, `lower`, and `initcap` functions
- Pattern matching operators (`LIKE`, `SIMILAR TO`, and POSIX-style regular expressions); locales affect both case insensitive matching and the classification of characters by character-class regular expressions
- The `to_char` family of functions
- The ability to use indexes with `LIKE` clauses

The drawback of using locales other than `C` or `POSIX` in Postgres Pro is its performance impact. It slows character handling and prevents ordinary indexes from being used by `LIKE`. For this reason use locales only if you actually need them.

As a workaround to allow Postgres Pro to use indexes with `LIKE` clauses under a non-`C` locale, several custom operator classes exist. These allow the creation of an index that performs a strict character-by-character comparison, ignoring locale comparison rules. Refer to [Section 11.10](#) for more information. Another approach is to create indexes using the `C` collation, as discussed in [Section 23.2](#).

23.1.3. Selecting Locales

Locales can be selected in different scopes depending on requirements. The above overview showed how locales are specified using `initdb` to set the defaults for the entire cluster. The following list shows where locales can be selected. Each item provides the defaults for the subsequent items, and each lower item allows overriding the defaults on a finer granularity.

1. As explained above, the environment of the operating system provides the defaults for the locales of a newly initialized database cluster. In many cases, this is enough: if the operating system is configured for the desired language/territory, by default Postgres Pro will also behave according to that locale.
2. As shown above, command-line options for `initdb` specify the locale settings for a newly initialized database cluster. Use this if the operating system does not have the locale configuration you want for your database system.
3. A locale can be selected separately for each database. The SQL command `CREATE DATABASE` and its command-line equivalent `createdb` have options for that. Use this for example if a database cluster houses databases for multiple tenants with different requirements.
4. Locale settings can be made for individual table columns. This uses an SQL object called *collation* and is explained in [Section 23.2](#). Use this for example to sort data in different languages or customize the sort order of a particular table.
5. Finally, locales can be selected for an individual query. Again, this uses SQL collation objects. This could be used to change the sort order based on run-time choices or for ad-hoc experimentation.

23.1.4. Locale Providers

A locale provider specifies which library defines the locale behavior for collations and character classifications.

The commands and tools that select the locale settings, as described above, each have an option to select the locale provider. Here is an example to initialize a database cluster using the ICU provider:

```
initdb --locale-provider=icu --icu-locale=en
```

See the description of the respective commands and programs for details. Note that you can mix locale providers at different granularities, for example use `libc` by default for the cluster but have one database that uses the `icu` provider, and then have collation objects using either provider within those databases.

Regardless of the locale provider, the operating system is still used to provide some locale-aware behavior, such as messages (see [lc_messages](#)).

The available locale providers are listed below:

`builtin`

The `builtin` provider uses built-in operations. Only the `C` and `C.UTF-8` locales are supported for this provider.

The `C` locale behavior is identical to the `C` locale in the `libc` provider. When using this locale, the behavior may depend on the database encoding.

The `C.UTF-8` locale is available only for when the database encoding is `UTF-8`, and the behavior is based on Unicode. The collation uses the code point values only. The regular expression character classes are based on the "POSIX Compatible" semantics, and the case mapping is the "simple" variant.

`icu`

The `icu` provider uses the external ICU library. Postgres Pro must have been configured with support.

ICU provides collation and character classification behavior that is independent of the operating system and database encoding, which is preferable if you expect to transition to other platforms without any change in results. `LC_COLLATE` and `LC_CTYPE` can be set independently of the ICU locale.

Note

For the ICU provider, results may depend on the version of the ICU library used, as it is updated to reflect changes in natural language over time.

`libc`

The `libc` provider uses the operating system's C library. The collation and character classification behavior is controlled by the settings `LC_COLLATE` and `LC_CTYPE`, so they cannot be set independently.

Note

The same locale name may have different behavior on different platforms when using the `libc` provider.

23.1.5. ICU Locales

23.1.5.1. ICU Locale Names

The ICU format for the locale name is a [Language Tag](#).

```
CREATE COLLATION mycollation1 (provider = icu, locale = 'ja-JP');
CREATE COLLATION mycollation2 (provider = icu, locale = 'fr');
```

23.1.5.2. Locale Canonicalization and Validation

When defining a new ICU collation object or database with ICU as the provider, the given locale name is transformed ("canonicalized") into a language tag if not already in that form. For instance,

```
CREATE COLLATION mycollation3 (provider = icu, locale = 'en-US-u-ku-tru');
NOTICE: using standard form "en-US-u-ku" for locale "en-US-u-ku-tru"
CREATE COLLATION mycollation4 (provider = icu, locale = 'de-DE.utf8');
NOTICE: using standard form "de-DE" for locale "de-DE.utf8"
```

If you see this notice, ensure that the `provider` and `locale` are the expected result. For consistent results when using the ICU provider, specify the canonical [language tag](#) instead of relying on the transformation.

A locale with no language name, or the special language name `root`, is transformed to have the language `und` ("undefined").

ICU can transform most `libc` locale names, as well as some other formats, into language tags for easier transition to ICU. If a `libc` locale name is used in ICU, it may not have precisely the same behavior as in `libc`.

If there is a problem interpreting the locale name, or if the locale name represents a language or region that ICU does not recognize, you will see the following warning:

```
CREATE COLLATION nonsense (provider = icu, locale = 'nonsense');
WARNING: ICU locale "nonsense" has unknown language "nonsense"
HINT: To disable ICU locale validation, set parameter icu_validation_level to
DISABLED.
CREATE COLLATION
```

[icu_validation_level](#) controls how the message is reported. Unless set to `ERROR`, the collation will still be created, but the behavior may not be what the user intended.

23.1.5.3. Language Tag

A language tag, defined in BCP 47, is a standardized identifier used to identify languages, regions, and other information about a locale.

Basic language tags are simply *language-region*; or even just *language*. The *language* is a language code (e.g. `fr` for French), and *region* is a region code (e.g. `CA` for Canada). Examples: `ja-JP`, `de`, or `fr-CA`.

Collation settings may be included in the language tag to customize collation behavior. ICU allows extensive customization, such as sensitivity (or insensitivity) to accents, case, and punctuation; treatment of digits within text; and many other options to satisfy a variety of uses.

To include this additional collation information in a language tag, append `-u`, which indicates there are additional collation settings, followed by one or more *-key-value* pairs. The *key* is the key for a [collation setting](#) and *value* is a valid value for that setting. For boolean settings, the *-key* may be specified without a corresponding *-value*, which implies a value of `true`.

For example, the language tag `en-US-u-kn-ks-level2` means the locale with the English language in the US region, with collation settings `kn` set to `true` and `ks` set to `level2`. Those settings mean the collation will be case-insensitive and treat a sequence of digits as a single number:

```
CREATE COLLATION mycollation5 (provider = icu, deterministic = false, locale = 'en-US-u-kn-ks-level2');
SELECT 'aB' = 'Ab' COLLATE mycollation5 as result;
   result
-----
      t
(1 row)

SELECT 'N-45' < 'N-123' COLLATE mycollation5 as result;
   result
-----
      t
(1 row)
```

See [Section 23.2.3](#) for details and additional examples of using language tags with custom collation information for the locale.

23.1.6. Problems

If locale support doesn't work according to the explanation above, check that the locale support in your operating system is correctly configured. To check what locales are installed on your system, you can use the command `locale -a` if your operating system provides it.

Check that Postgres Pro is actually using the locale that you think it is. The `LC_COLLATE` and `LC_CTYPE` settings are determined when a database is created, and cannot be changed except by creating a new database. Other locale settings including `LC_MESSAGES` and `LC_MONETARY` are initially determined by the environment the server is started in, but can be changed on-the-fly. You can check the active locale settings using the `SHOW` command.

Client applications that handle server-side errors by parsing the text of the error message will obviously have problems when the server's messages are in a different language. Authors of such applications are advised to make use of the error code scheme instead.

23.2. Collation Support

The collation feature allows specifying the sort order and character classification behavior of data per-column, or even per-operation. This alleviates the restriction that the `LC_COLLATE` and `LC_CTYPE` settings of a database cannot be changed after its creation.

23.2.1. Concepts

Conceptually, every expression of a collatable data type has a collation. (The built-in collatable data types are `text`, `varchar`, and `char`. User-defined base types can also be marked collatable, and of course a [domain](#) over a collatable data type is collatable.) If the expression is a column reference, the collation

of the expression is the defined collation of the column. If the expression is a constant, the collation is the default collation of the data type of the constant. The collation of a more complex expression is derived from the collations of its inputs, as described below.

The collation of an expression can be the “default” collation, which means the locale settings defined for the database. It is also possible for an expression's collation to be indeterminate. In such cases, ordering operations and other operations that need to know the collation will fail.

When the database system has to perform an ordering or a character classification, it uses the collation of the input expression. This happens, for example, with `ORDER BY` clauses and function or operator calls such as `<`. The collation to apply for an `ORDER BY` clause is simply the collation of the sort key. The collation to apply for a function or operator call is derived from the arguments, as described below. In addition to comparison operators, collations are taken into account by functions that convert between lower and upper case letters, such as `lower`, `upper`, and `initcap`; by pattern matching operators; and by `to_char` and related functions.

For a function or operator call, the collation that is derived by examining the argument collations is used at run time for performing the specified operation. If the result of the function or operator call is of a collatable data type, the collation is also used at parse time as the defined collation of the function or operator expression, in case there is a surrounding expression that requires knowledge of its collation.

The *collation derivation* of an expression can be implicit or explicit. This distinction affects how collations are combined when multiple different collations appear in an expression. An explicit collation derivation occurs when a `COLLATE` clause is used; all other collation derivations are implicit. When multiple collations need to be combined, for example in a function call, the following rules are used:

1. If any input expression has an explicit collation derivation, then all explicitly derived collations among the input expressions must be the same, otherwise an error is raised. If any explicitly derived collation is present, that is the result of the collation combination.
2. Otherwise, all input expressions must have the same implicit collation derivation or the default collation. If any non-default collation is present, that is the result of the collation combination. Otherwise, the result is the default collation.
3. If there are conflicting non-default implicit collations among the input expressions, then the combination is deemed to have indeterminate collation. This is not an error condition unless the particular function being invoked requires knowledge of the collation it should apply. If it does, an error will be raised at run-time.

For example, consider this table definition:

```
CREATE TABLE test1 (
    a text COLLATE "ru_RU",
    b text COLLATE "es_ES",
    ...
);
```

Then in

```
SELECT a < 'foo' FROM test1;
```

the `<` comparison is performed according to `ru_RU` rules, because the expression combines an implicitly derived collation with the default collation. But in

```
SELECT a < ('foo' COLLATE "fr_FR") FROM test1;
```

the comparison is performed using `fr_FR` rules, because the explicit collation derivation overrides the implicit one. Furthermore, given

```
SELECT a < b FROM test1;
```

the parser cannot determine which collation to apply, since the `a` and `b` columns have conflicting implicit collations. Since the `<` operator does need to know which collation to use, this will result in an error. The error can be resolved by attaching an explicit collation specifier to either input expression, thus:

```
SELECT a < b COLLATE "ru_RU" FROM test1;
```

or equivalently

```
SELECT a COLLATE "ru_RU" < b FROM test1;
```

On the other hand, the structurally similar case

```
SELECT a || b FROM test1;
```

does not result in an error, because the `||` operator does not care about collations: its result is the same regardless of the collation.

The collation assigned to a function or operator's combined input expressions is also considered to apply to the function or operator's result, if the function or operator delivers a result of a collatable data type. So, in

```
SELECT * FROM test1 ORDER BY a || 'foo';
```

the ordering will be done according to `ru_RU` rules. But this query:

```
SELECT * FROM test1 ORDER BY a || b;
```

results in an error, because even though the `||` operator doesn't need to know a collation, the `ORDER BY` clause does. As before, the conflict can be resolved with an explicit collation specifier:

```
SELECT * FROM test1 ORDER BY a || b COLLATE "fr_FR";
```

23.2.2. Managing Collations

A collation is an SQL schema object that maps an SQL name to locales provided by libraries installed in the operating system. A collation definition has a *provider* that specifies which library supplies the locale data. One standard provider name is `libc`, which uses the locales provided by the operating system C library. These are the locales used by most tools provided by the operating system. Another provider is `icu`, which uses the external ICU library. ICU locales can only be used if support for ICU was configured when Postgres Pro was built.

A collation object provided by `libc` maps to a combination of `LC_COLLATE` and `LC_CTYPE` settings, as accepted by the `setlocale()` system library call. (As the name would suggest, the main purpose of a collation is to set `LC_COLLATE`, which controls the sort order. But it is rarely necessary in practice to have an `LC_CTYPE` setting that is different from `LC_COLLATE`, so it is more convenient to collect these under one concept than to create another infrastructure for setting `LC_CTYPE` per expression.) Also, a `libc` collation is tied to a character set encoding (see [Section 23.3](#)). The same collation name may exist for different encodings.

A collation object provided by `icu` maps to a named collator provided by the ICU library. ICU does not support separate “collate” and “ctype” settings, so they are always the same. Also, ICU collations are independent of the encoding, so there is always only one ICU collation of a given name in a database.

23.2.2.1. Standard Collations

On all platforms, the following collations are supported:

`unicode`

This SQL standard collation sorts using the Unicode Collation Algorithm with the Default Unicode Collation Element Table. It is available in all encodings. ICU support is required to use this collation, and behavior may change if Postgres Pro is built with a different version of ICU. (This collation has the same behavior as the ICU root locale; see [und-x-icu](#) (for “undefined”).)

`ucs_basic`

This SQL standard collation sorts using the Unicode code point values rather than natural language order, and only the ASCII letters “A” through “Z” are treated as letters. The behavior is efficient and

stable across all versions. Only available for encoding `UTF8`. (This collation has the same behavior as the `libc` locale specification `C` in `UTF8` encoding.)

`pg_c_utf8`

This collation sorts by Unicode code point values rather than natural language order. For the functions `lower`, `initcap`, and `upper`, it uses Unicode simple case mapping. For pattern matching (including regular expressions), it uses the POSIX Compatible variant of Unicode [Compatibility Properties](#). Behavior is efficient and stable within a Postgres Pro major version. This collation is only available for encoding `UTF8`.

`C` (equivalent to `POSIX`)

The `C` and `POSIX` collations are based on “traditional C” behavior. They sort by byte values rather than natural language order, and only the ASCII letters “A” through “Z” are treated as letters. The behavior is efficient and stable across all versions for a given database encoding, but behavior may vary between different database encodings.

`default`

The `default` collation selects the locale specified at database creation time.

Additional collations may be available depending on operating system support. The efficiency and stability of these additional collations depend on the collation provider, the provider version, and the locale.

23.2.2.2. Predefined Collations

If the operating system provides support for using multiple locales within a single program (`newlocale` and related functions), or if support for ICU is configured, then when a database cluster is initialized, `initdb` populates the system catalog `pg_collation` with collations based on all the locales it finds in the operating system at the time.

To inspect the currently available locales, use the query `SELECT * FROM pg_collation`, or the command `\dos+ in psql`.

23.2.2.2.1. libc Collations

For example, the operating system might provide a locale named `ru_RU.utf8`. `initdb` would then create a collation named `ru_RU.utf8` for encoding `UTF8` that has both `LC_COLLATE` and `LC_CTYPE` set to `ru_RU.utf8`. It will also create a collation with the `.utf8` tag stripped off the name. So you could also use the collation under the name `ru_RU`, which is less cumbersome to write and makes the name less encoding-dependent. Note that, nevertheless, the initial set of collation names is platform-dependent.

The default set of collations provided by `libc` map directly to the locales installed in the operating system, which can be listed using the command `locale -a`. In case a `libc` collation is needed that has different values for `LC_COLLATE` and `LC_CTYPE`, or if new locales are installed in the operating system after the database system was initialized, then a new collation may be created using the [CREATE COLLATION](#) command. New operating system locales can also be imported en masse using the [pg_import_system_collations\(\)](#) function.

Within any particular database, only collations that use that database's encoding are of interest. Other entries in `pg_collation` are ignored. Thus, a stripped collation name such as `ru_RU` can be considered unique within a given database even though it would not be unique globally. Use of the stripped collation names is recommended, since it will make one fewer thing you need to change if you decide to change to another database encoding. Note however that the `default`, `C`, and `POSIX` collations can be used regardless of the database encoding.

Postgres Pro considers distinct collation objects to be incompatible even when they have identical properties. Thus for example,

```
SELECT a COLLATE "C" < b COLLATE "POSIX" FROM test1;
```

will draw an error even though the `C` and `POSIX` collations have identical behaviors. Mixing stripped and non-stripped collation names is therefore not recommended.

23.2.2.2.2. ICU Collations

With ICU, it is not sensible to enumerate all possible locale names. ICU uses a particular naming system for locales, but there are many more ways to name a locale than there are actually distinct locales. `initdb` uses the ICU APIs to extract a set of distinct locales to populate the initial set of collations. Collations provided by ICU are created in the SQL environment with names in BCP 47 language tag format, with a “private use” extension `-x-icu` appended, to distinguish them from `libc` locales.

Here are some example collations that might be created:

`ru-x-icu`

Russian collation, default variant

`ru-UA-x-icu`

Russian collation for Ukraine, default variant

(There is also, say, `ru-RU-x-icu`, but as of this writing, it is equivalent to `ru-x-icu`.)

`und-x-icu` (for “undefined”)

ICU “root” collation. Use this to get a reasonable language-agnostic sort order.

Some (less frequently used) encodings are not supported by ICU. When the database encoding is one of these, ICU collation entries in `pg_collation` are ignored. Attempting to use one will draw an error along the lines of “collation “de-x-icu” for encoding “WIN874” does not exist”.

23.2.2.3. Creating New Collation Objects

If the standard and predefined collations are not sufficient, users can create their own collation objects using the SQL command [CREATE COLLATION](#).

The standard and predefined collations are in the schema `pg_catalog`, like all predefined objects. User-defined collations should be created in user schemas. This also ensures that they are saved by `pg_dump`.

23.2.2.3.1. libc Collations

New `libc` collations can be created like this:

```
CREATE COLLATION russian (provider = libc, locale = 'ru_RU');
```

The exact values that are acceptable for the `locale` clause in this command depend on the operating system. On Unix-like systems, the command `locale -a` will show a list.

Since the predefined `libc` collations already include all collations defined in the operating system when the database instance is initialized, it is not often necessary to manually create new ones. Reasons might be if a different naming system is desired (in which case see also [Section 23.2.2.3.3](#)) or if the operating system has been upgraded to provide new locale definitions (in which case see also `pg_import_system_collations()`).

23.2.2.3.2. ICU Collations

ICU collations can be created like:

```
CREATE COLLATION german (provider = icu, locale = 'de-DE');
```

ICU locales are specified as a BCP 47 [Language Tag](#), but can also accept most `libc`-style locale names. If possible, `libc`-style locale names are transformed into language tags.

New ICU collations can customize collation behavior extensively by including collation attributes in the language tag. See [Section 23.2.3](#) for details and examples.

23.2.2.3.3. Copying Collations

The command `CREATE COLLATION` can also be used to create a new collation from an existing collation, which can be useful to be able to use operating-system-independent collation names in applications, create compatibility names, or use an ICU-provided collation under a more readable name. For example:

```
CREATE COLLATION russian FROM "ru_RU";
CREATE COLLATION french FROM "fr-x-icu";
```

23.2.2.4. Nondeterministic Collations

A collation is either *deterministic* or *nondeterministic*. A deterministic collation uses deterministic comparisons, which means that it considers strings to be equal only if they consist of the same byte sequence. Nondeterministic comparison may determine strings to be equal even if they consist of different bytes. Typical situations include case-insensitive comparison, accent-insensitive comparison, as well as comparison of strings in different Unicode normal forms. It is up to the collation provider to actually implement such insensitive comparisons; the deterministic flag only determines whether ties are to be broken using bytewise comparison. See also [Unicode Technical Standard 10](#) for more information on the terminology.

To create a nondeterministic collation, specify the property `deterministic = false` to `CREATE COLLATION`, for example:

```
CREATE COLLATION ndcoll (provider = icu, locale = 'und', deterministic = false);
```

This example would use the standard Unicode collation in a nondeterministic way. In particular, this would allow strings in different normal forms to be compared correctly. More interesting examples make use of the ICU customization facilities explained above. For example:

```
CREATE COLLATION case_insensitive (provider = icu, locale = 'und-u-ks-level2',
    deterministic = false);
CREATE COLLATION ignore_accents (provider = icu, locale = 'und-u-ks-level1-kc-true',
    deterministic = false);
```

All standard and predefined collations are deterministic, all user-defined collations are deterministic by default. While nondeterministic collations give a more “correct” behavior, especially when considering the full power of Unicode and its many special cases, they also have some drawbacks. Foremost, their use leads to a performance penalty. Note, in particular, that B-tree cannot use deduplication with indexes that use a nondeterministic collation. Also, certain operations are not possible with nondeterministic collations, such as pattern matching operations. Therefore, they should be used only in cases where they are specifically wanted.

Tip

To deal with text in different Unicode normalization forms, it is also an option to use the functions/expressions `normalize` and `is_normalized` to preprocess or check the strings, instead of using nondeterministic collations. There are different trade-offs for each approach.

23.2.3. ICU Custom Collations

ICU allows extensive control over collation behavior by defining new collations with collation settings as a part of the language tag. These settings can modify the collation order to suit a variety of needs. For instance:

```
-- ignore differences in accents and case
CREATE COLLATION ignore_accent_case (provider = icu, deterministic = false, locale =
    'und-u-ks-level1');
```



```

SELECT 'Å' = 'A' COLLATE ignore_accent_case; -- true
SELECT 'z' = 'Z' COLLATE ignore_accent_case; -- true

-- upper case letters sort before lower case.
CREATE COLLATION upper_first (provider = icu, locale = 'und-u-kf-upper');
SELECT 'B' < 'b' COLLATE upper_first; -- true

-- treat digits numerically and ignore punctuation
CREATE COLLATION num_ignore_punct (provider = icu, deterministic = false, locale =
    'und-u-ka-shifted-kn');
SELECT 'id-45' < 'id-123' COLLATE num_ignore_punct; -- true
SELECT 'w;x*y-z' = 'wxyz' COLLATE num_ignore_punct; -- true

```

Many of the available options are described in [Section 23.2.3.2](#), or see [Section 23.2.3.5](#) for more details.

23.2.3.1. ICU Comparison Levels

Comparison of two strings (collation) in ICU is determined by a multi-level process, where textual features are grouped into "levels". Treatment of each level is controlled by the [collation settings](#). Higher levels correspond to finer textual features.

[Table 23.1](#) shows which textual feature differences are considered significant when determining equality at the given level. The Unicode character U+2063 is an invisible separator, and as seen in the table, is ignored for at all levels of comparison less than `identic`.

Table 23.1. ICU Collation Levels

Level	Description	'f' = 'f'	'ab' = U&'a\2063b'	'x-y' = 'x_y'	'g' = 'G'	'n' = 'ñ'	'y' = 'z'
level1	Base Character	true	true	true	true	true	false
level2	Accents	true	true	true	true	false	false
level3	Case/Variants	true	true	true	false	false	false
level4	Punctuation ^a	true	true	false	false	false	false
identic	All	true	false	false	false	false	false

^aonly with `ka-shifted`; see [Table 23.2](#)

At every level, even with full normalization off, basic normalization is performed. For example, 'á' may be composed of the code points U&'0061\0301' or the single code point U&'00E1', and those sequences will be considered equal even at the `identic` level. To treat any difference in code point representation as distinct, use a collation created with `deterministic` set to `true`.

23.2.3.1.1. Collation Level Examples

```

CREATE COLLATION level3 (provider = icu, deterministic = false, locale = 'und-u-ka-shifted-ks-level3');
CREATE COLLATION level4 (provider = icu, deterministic = false, locale = 'und-u-ka-shifted-ks-level4');
CREATE COLLATION identic (provider = icu, deterministic = false, locale = 'und-u-ka-shifted-ks-identic');

-- invisible separator ignored at all levels except identic
SELECT 'ab' = U&'a\2063b' COLLATE level4; -- true
SELECT 'ab' = U&'a\2063b' COLLATE identic; -- false

-- punctuation ignored at level3 but not at level 4

```

```
SELECT 'x-y' = 'x_y' COLLATE level3; -- true
SELECT 'x-y' = 'x_y' COLLATE level4; -- false
```

23.2.3.2. Collation Settings for an ICU Locale

[Table 23.2](#) shows the available collation settings, which can be used as part of a language tag to customize a collation.

Table 23.2. ICU Collation Settings

Key	Values	Default	Description
co	emoji, phonebk, standard, ...	standard	Collation type. See Section 23.2.3.5 for additional options and details.
ka	noignore, shifted	noignore	If set to <code>shifted</code> , causes some characters (e.g. punctuation or space) to be ignored in comparison. Key <code>ks</code> must be set to <code>level3</code> or lower to take effect. Set key <code>kv</code> to control which character classes are ignored.
kb	true, false	false	Backwards comparison for the level 2 differences. For example, locale <code>und-u-kb</code> sorts 'àe' before 'aé'.
kc	true, false	false	Separates case into a "level 2.5" that falls between accents and other level 3 features. If set to <code>true</code> and <code>ks</code> is set to <code>level1</code> , will ignore accents but take case into account.
kf	upper, lower, false	false	If set to <code>upper</code> , upper case sorts before lower case. If set to <code>lower</code> , lower case sorts before upper case. If set to <code>false</code> , the sort depends on the rules of the locale.
kn	true, false	false	If set to <code>true</code> , numbers within a string are treated as a single numeric value rather than a sequence of digits. For example, 'id-45' sorts before 'id-123'.
kk	true, false	false	Enable full normalization; may affect performance. Basic normalization is performed even when set to <code>false</code> . Locales for languages that require full normalization typically enable it by default. Full normalization is important in some cases, such as when multiple accents are applied to a single character. For example, the code point sequences <code>U&'\0065\0323\0302'</code> and <code>U&'\0065\0302\0323'</code> represent an <code>e</code> with circumflex and dot-below accents applied in different orders. With full normalization on, these code point sequences are treated as equal; otherwise they are unequal.
kr	space, punct, symbol, currency, digit, <i>script-id</i>		Set to one or more of the valid values, or any BCP 47 <i>script-id</i> , e.g. <code>latn</code> ("Latin") or <code>grek</code> ("Greek"). Multiple values are separated by "-". Redefines the ordering of classes of characters; those characters belonging to a class earlier in the list sort before characters belonging to a

Key	Values	Default	Description
			class later in the list. For instance, the value <code>digit-currency-space</code> (as part of a language tag like <code>und-u-kr-digit-currency-space</code>) sorts punctuation before digits and spaces.
<code>ks</code>	<code>level1</code> , <code>level2</code> , <code>level3</code> , <code>level4</code> , <code>identic</code>	<code>level3</code>	Sensitivity (or "strength") when determining equality, with <code>level1</code> the least sensitive to differences and <code>identic</code> the most sensitive to differences. See Table 23.1 for details.
<code>kv</code>	<code>space</code> , <code>punct</code> , <code>symbol</code> , <code>currency</code>	<code>punct</code>	Classes of characters ignored during comparison at level 3. Setting to a later value includes earlier values; e.g. <code>symbol</code> also includes <code>punct</code> and <code>space</code> in the characters to be ignored. Key <code>ka</code> must be set to <code>shifted</code> and key <code>ks</code> must be set to <code>level3</code> or lower to take effect.

Defaults may depend on locale. The above table is not meant to be complete. See [Section 23.2.3.5](#) for additional options and details.

Note

For many collation settings, you must create the collation with `deterministic` set to `false` for the setting to have the desired effect (see [Section 23.2.2.4](#)). Additionally, some settings only take effect when the key `ka` is set to `shifted` (see [Table 23.2](#)).

23.2.3.3. Collation Settings Examples

```
CREATE COLLATION "de-u-co-phonebk-x-icu" (provider = icu, locale = 'de-u-co-phonebk');
```

German collation with phone book collation type

```
CREATE COLLATION "und-u-co-emoji-x-icu" (provider = icu, locale = 'und-u-co-emoji');
```

Root collation with Emoji collation type, per Unicode Technical Standard #51

```
CREATE COLLATION latn_cyrl (provider = icu, locale = 'ru-RU-u-kr-latn-cyrl');
```

Sort Latin letters before Cyrillic ones. (The default is Cyrillic before Latin.)

```
CREATE COLLATION upperfirst (provider = icu, locale = 'ru-RU-u-kf-upper');
```

Sort upper-case letters before lower-case letters. (The default is lower-case letters first.)

```
CREATE COLLATION special (provider = icu, locale = 'ru-RU-u-kf-upper-kr-latn-cyrl');
```

Combines both of the above options.

23.2.3.4. ICU Tailoring Rules

If the options provided by the collation settings shown above are not sufficient, the order of collation elements can be changed with tailoring rules, whose syntax is detailed at <https://unicode-org.github.io/icu/userguide/collation/customization/>.

This small example creates a collation based on the root locale with a tailoring rule:

```
CREATE COLLATION custom (provider = icu, locale = 'und', rules = '&V << w <<< W');
```

With this rule, the letter “W” is sorted after “V”, but is treated as a secondary difference similar to an accent. Rules like this are contained in the locale definitions of some languages. (Of course, if a locale definition already contains the desired rules, then they don't need to be specified again explicitly.)

Here is a more complex example. The following statement sets up a collation named `ebcdic` with rules to sort US-ASCII characters in the order of the EBCDIC encoding.

```
CREATE COLLATION ebcdic (provider = icu, locale = 'und',
rules = $$
& ' ' < '.' < '<' < '(' < '+' < \ |
< '&' < '!' < '$' < '*' < ')' < ';'
< '-' < '/' < ',' < '%' < '_' < '>' < '?'
< ':' < '#' < '@' < '\' < '=' < '"'
< *a-r < '~' < *s-z < '^' < '[' < ']'
< '{' < *A-I < '}' < *J-R < '\' < *S-Z < *0-9
$$);

SELECT c
FROM (VALUES ('a'), ('b'), ('A'), ('B'), ('1'), ('2'), ('!'), ('^')) AS x(c)
ORDER BY c COLLATE ebcdic;
 c
---
 !
 a
 b
 ^
 A
 B
 1
 2
```

23.2.3.5. External References for ICU

This section (Section 23.2.3) is only a brief overview of ICU behavior and language tags. Refer to the following documents for technical details, additional options, and new behavior:

- [Unicode Technical Standard #35](#)
- [BCP 47](#)
- [CLDR repository](#)
- <https://unicode-org.github.io/icu/userguide/locale/>
- <https://unicode-org.github.io/icu/userguide/collation/>

23.3. Character Set Support

The character set support in Postgres Pro allows you to store text in a variety of character sets (also called encodings), including single-byte character sets such as the ISO 8859 series and multiple-byte character sets such as EUC (Extended Unix Code), UTF-8, and Mule internal code. All supported character sets can be used transparently by clients, but a few are not supported for use within the server (that is, as a server-side encoding). The default character set is selected while initializing your Postgres Pro database cluster using `initdb`. It can be overridden when you create a database, so you can have multiple databases each with a different character set.

An important restriction, however, is that each database's character set must be compatible with the database's `LC_CTYPE` (character classification) and `LC_COLLATE` (string sort order) locale settings. For `C` or `POSIX` locale, any character set is allowed, but for other libc-provided locales there is only one character set that will work correctly. (On Windows, however, UTF-8 encoding can be used with any

locale.) If you have ICU support configured, ICU-provided locales can be used with most but not all server-side encodings.

23.3.1. Supported Character Sets

[Table 23.3](#) shows the character sets available for use in Postgres Pro.

Table 23.3. Postgres Pro Character Sets

Name	Description	Language	Server?	ICU?	Bytes/ Char	Aliases
BIG5	Big Five	Traditional Chinese	No	No	1-2	WIN950, Windows950
EUC_CN	Extended UNIX Code-CN	Simplified Chinese	Yes	Yes	1-3	
EUC_JP	Extended UNIX Code-JP	Japanese	Yes	Yes	1-3	
EUC_JIS_2004	Extended UNIX Code-JP, JIS X 0213	Japanese	Yes	No	1-3	
EUC_KR	Extended UNIX Code-KR	Korean	Yes	Yes	1-3	
EUC_TW	Extended UNIX Code-TW	Traditional Chinese, Taiwanese	Yes	Yes	1-4	
GB18030	National Standard	Chinese	No	No	1-4	
GBK	Extended National Standard	Simplified Chinese	No	No	1-2	WIN936, Windows936
ISO_8859_5	ISO 8859-5, ECMA 113	Latin/Cyrillic	Yes	Yes	1	
ISO_8859_6	ISO 8859-6, ECMA 114	Latin/Arabic	Yes	Yes	1	
ISO_8859_7	ISO 8859-7, ECMA 118	Latin/Greek	Yes	Yes	1	
ISO_8859_8	ISO 8859-8, ECMA 121	Latin/Hebrew	Yes	Yes	1	
JOHAB	JOHAB	Korean (Hangul)	No	No	1-3	
KOI8R	KOI8-R	Cyrillic (Russian)	Yes	Yes	1	KOI8
KOI8U	KOI8-U	Cyrillic (Ukrainian)	Yes	Yes	1	
LATIN1	ISO 8859-1, ECMA 94	Western European	Yes	Yes	1	ISO88591
LATIN2	ISO 8859-2, ECMA 94	Central European	Yes	Yes	1	ISO88592
LATIN3	ISO 8859-3, ECMA 94	South European	Yes	Yes	1	ISO88593
LATIN4	ISO 8859-4, ECMA 94	North European	Yes	Yes	1	ISO88594

Name	Description	Language	Server?	ICU?	Bytes/ Char	Aliases
LATIN5	ISO 8859-9, ECMA 128	Turkish	Yes	Yes	1	ISO88599
LATIN6	ISO 8859-10, ECMA 144	Nordic	Yes	Yes	1	ISO885910
LATIN7	ISO 8859-13	Baltic	Yes	Yes	1	ISO885913
LATIN8	ISO 8859-14	Celtic	Yes	Yes	1	ISO885914
LATIN9	ISO 8859-15	LATIN1 with Euro and ac- cents	Yes	Yes	1	ISO885915
LATIN10	ISO 8859-16, ASRO SR 14111	Romanian	Yes	No	1	ISO885916
MULE_INTERNAL	Mule internal code	Multilingual Emacs	Yes	No	1-4	
SJIS	Shift JIS	Japanese	No	No	1-2	Mskanji, ShiftJIS, WIN932, Win- dows932
SHIFT_JIS_2004	Shift JIS, JIS X 0213	Japanese	No	No	1-2	
SQL_ASCII	unspecified (see text)	<i>any</i>	Yes	No	1	
UHC	Unified Hangul Code	Korean	No	No	1-2	WIN949, Win- dows949
UTF8	Unicode, 8-bit	<i>all</i>	Yes	Yes	1-4	Unicode
WIN866	Windows CP866	Cyrillic	Yes	Yes	1	ALT
WIN874	Windows CP874	Thai	Yes	No	1	
WIN1250	Windows CP1250	Central Euro- pean	Yes	Yes	1	
WIN1251	Windows CP1251	Cyrillic	Yes	Yes	1	WIN
WIN1252	Windows CP1252	Western Euro- pean	Yes	Yes	1	
WIN1253	Windows CP1253	Greek	Yes	Yes	1	
WIN1254	Windows CP1254	Turkish	Yes	Yes	1	
WIN1255	Windows CP1255	Hebrew	Yes	Yes	1	
WIN1256	Windows CP1256	Arabic	Yes	Yes	1	
WIN1257	Windows CP1257	Baltic	Yes	Yes	1	

Name	Description	Language	Server?	ICU?	Bytes/ Char	Aliases
WIN1258	Windows CP1258	Vietnamese	Yes	Yes	1	ABC, TCVN, TCVN5712, VSCII

Not all client APIs support all the listed character sets. For example, the Postgres Pro JDBC driver does not support `MULE_INTERNAL`, `LATIN6`, `LATIN8`, and `LATIN10`.

The `SQL_ASCII` setting behaves considerably differently from the other settings. When the server character set is `SQL_ASCII`, the server interprets byte values 0–127 according to the ASCII standard, while byte values 128–255 are taken as uninterpreted characters. No encoding conversion will be done when the setting is `SQL_ASCII`. Thus, this setting is not so much a declaration that a specific encoding is in use, as a declaration of ignorance about the encoding. In most cases, if you are working with any non-ASCII data, it is unwise to use the `SQL_ASCII` setting because Postgres Pro will be unable to help you by converting or validating non-ASCII characters.

23.3.2. Setting the Character Set

`initdb` defines the default character set (encoding) for a Postgres Pro cluster. For example,

```
initdb -E EUC_JP
```

sets the default character set to `EUC_JP` (Extended Unix Code for Japanese). You can use `--encoding` instead of `-E` if you prefer longer option strings. If no `-E` or `--encoding` option is given, `initdb` attempts to determine the appropriate encoding to use based on the specified or default locale.

You can specify a non-default encoding at database creation time, provided that the encoding is compatible with the selected locale:

```
createdb -E EUC_KR -T template0 --lc-collate=ko_KR.euckr --lc-ctype=ko_KR.euckr korean
```

This will create a database named `korean` that uses the character set `EUC_KR`, and locale `ko_KR`. Another way to accomplish this is to use this SQL command:

```
CREATE DATABASE korean WITH ENCODING 'EUC_KR' LC_COLLATE='ko_KR.euckr'
LC_CTYPE='ko_KR.euckr' TEMPLATE=template0;
```

Notice that the above commands specify copying the `template0` database. When copying any other database, the encoding and locale settings cannot be changed from those of the source database, because that might result in corrupt data. For more information see [Section 22.4](#).

The encoding for a database is stored in the system catalog `pg_database`. You can see it by using the `psql -l` option or the `\l` command.

```
$ psql -l
```

```

                                List of databases
  Name      | Owner   | Encoding | Collation | Ctype   | Access
Privileges
-----+-----+-----+-----+-----+-----
 clocaledb | hlinnaka | SQL_ASCII | C         | C       |
 englishdb | hlinnaka | UTF8      | en_GB.UTF8 | en_GB.UTF8 |
 japanese  | hlinnaka | UTF8      | ja_JP.UTF8 | ja_JP.UTF8 |
 korean     | hlinnaka | EUC_KR    | ko_KR.euckr | ko_KR.euckr |
 postgres  | hlinnaka | UTF8      | fi_FI.UTF8 | fi_FI.UTF8 |
 template0 | hlinnaka | UTF8      | fi_FI.UTF8 | fi_FI.UTF8 | {=c/
hlinnaka,hlinnaka=CTc/hlinnaka}
 template1 | hlinnaka | UTF8      | fi_FI.UTF8 | fi_FI.UTF8 | {=c/
hlinnaka,hlinnaka=CTc/hlinnaka}
(7 rows)
```

Important

On most modern operating systems, Postgres Pro can determine which character set is implied by the `LC_CTYPE` setting, and it will enforce that only the matching database encoding is used. On older systems it is your responsibility to ensure that you use the encoding expected by the locale you have selected. A mistake in this area is likely to lead to strange behavior of locale-dependent operations such as sorting.

Postgres Pro will allow superusers to create databases with `SQL_ASCII` encoding even when `LC_CTYPE` is not `C` or `POSIX`. As noted above, `SQL_ASCII` does not enforce that the data stored in the database has any particular encoding, and so this choice poses risks of locale-dependent misbehavior. Using this combination of settings is deprecated and may someday be forbidden altogether.

23.3.3. Automatic Character Set Conversion Between Server and Client

Postgres Pro supports automatic character set conversion between server and client for many combinations of character sets ([Section 23.3.4](#) shows which ones).

To enable automatic character set conversion, you have to tell Postgres Pro the character set (encoding) you would like to use in the client. There are several ways to accomplish this:

- Using the `\encoding` command in `psql`. `\encoding` allows you to change client encoding on the fly. For example, to change the encoding to `SJIS`, type:

```
\encoding SJIS
```

- `libpq` ([Section 35.11](#)) has functions to control the client encoding.
- Using `SET client_encoding TO`. Setting the client encoding can be done with this SQL command:

```
SET CLIENT_ENCODING TO 'value';
```

Also you can use the standard SQL syntax `SET NAMES` for this purpose:

```
SET NAMES 'value';
```

To query the current client encoding:

```
SHOW client_encoding;
```

To return to the default encoding:

```
RESET client_encoding;
```

- Using `PGCLIENTENCODING`. If the environment variable `PGCLIENTENCODING` is defined in the client's environment, that client encoding is automatically selected when a connection to the server is made. (This can subsequently be overridden using any of the other methods mentioned above.)
- Using the configuration variable `client_encoding`. If the `client_encoding` variable is set, that client encoding is automatically selected when a connection to the server is made. (This can subsequently be overridden using any of the other methods mentioned above.)

If the conversion of a particular character is not possible — suppose you chose `EUC_JP` for the server and `LATIN1` for the client, and some Japanese characters are returned that do not have a representation in `LATIN1` — an error is reported.

If the client character set is defined as `SQL_ASCII`, encoding conversion is disabled, regardless of the server's character set. (However, if the server's character set is not `SQL_ASCII`, the server will still check that incoming data is valid for that encoding; so the net effect is as though the client character set were the same as the server's.) Just as for the server, use of `SQL_ASCII` is unwise unless you are working with all-ASCII data.

23.3.4. Available Character Set Conversions

Postgres Pro allows conversion between any two character sets for which a conversion function is listed in the `pg_conversion` system catalog. Postgres Pro comes with some predefined conversions, as summarized in [Table 23.4](#) and shown in more detail in [Table 23.5](#). You can create a new conversion using the SQL command `CREATE CONVERSION`. (To be used for automatic client/server conversions, a conversion must be marked as “default” for its character set pair.)

Table 23.4. Built-in Client/Server Character Set Conversions

Server Character Set	Available Client Character Sets
BIG5	<i>not supported as a server encoding</i>
EUC_CN	<i>EUC_CN</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
EUC_JP	<i>EUC_JP</i> , <i>MULE_INTERNAL</i> , <i>SJIS</i> , <i>UTF8</i>
EUC_JIS_2004	<i>EUC_JIS_2004</i> , <i>SHIFT_JIS_2004</i> , <i>UTF8</i>
EUC_KR	<i>EUC_KR</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
EUC_TW	<i>EUC_TW</i> , <i>BIG5</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
GB18030	<i>not supported as a server encoding</i>
GBK	<i>not supported as a server encoding</i>
ISO_8859_5	<i>ISO_8859_5</i> , <i>KOI8R</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i> , <i>WIN866</i> , <i>WIN1251</i>
ISO_8859_6	<i>ISO_8859_6</i> , <i>UTF8</i>
ISO_8859_7	<i>ISO_8859_7</i> , <i>UTF8</i>
ISO_8859_8	<i>ISO_8859_8</i> , <i>UTF8</i>
JOHAB	<i>not supported as a server encoding</i>
KOI8R	<i>KOI8R</i> , <i>ISO_8859_5</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i> , <i>WIN866</i> , <i>WIN1251</i>
KOI8U	<i>KOI8U</i> , <i>UTF8</i>
LATIN1	<i>LATIN1</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
LATIN2	<i>LATIN2</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i> , <i>WIN1250</i>
LATIN3	<i>LATIN3</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
LATIN4	<i>LATIN4</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i>
LATIN5	<i>LATIN5</i> , <i>UTF8</i>
LATIN6	<i>LATIN6</i> , <i>UTF8</i>
LATIN7	<i>LATIN7</i> , <i>UTF8</i>
LATIN8	<i>LATIN8</i> , <i>UTF8</i>
LATIN9	<i>LATIN9</i> , <i>UTF8</i>
LATIN10	<i>LATIN10</i> , <i>UTF8</i>
MULE_INTERNAL	<i>MULE_INTERNAL</i> , <i>BIG5</i> , <i>EUC_CN</i> , <i>EUC_JP</i> , <i>EUC_KR</i> , <i>EUC_TW</i> , <i>ISO_8859_5</i> , <i>KOI8R</i> , <i>LATIN1 to LATIN4</i> , <i>SJIS</i> , <i>WIN866</i> , <i>WIN1250</i> , <i>WIN1251</i>
SJIS	<i>not supported as a server encoding</i>
SHIFT_JIS_2004	<i>not supported as a server encoding</i>
SQL_ASCII	<i>any (no conversion will be performed)</i>
UHC	<i>not supported as a server encoding</i>
UTF8	<i>all supported encodings</i>
WIN866	<i>WIN866</i> , <i>ISO_8859_5</i> , <i>KOI8R</i> , <i>MULE_INTERNAL</i> , <i>UTF8</i> , <i>WIN1251</i>
WIN874	<i>WIN874</i> , <i>UTF8</i>

Server Character Set	Available Client Character Sets
WIN1250	<i>WIN1250</i> , LATIN2, MULE_INTERNAL , UTF8
WIN1251	<i>WIN1251</i> , ISO_8859_5 , KOI8R, MULE_INTERNAL , UTF8, WIN866
WIN1252	<i>WIN1252</i> , UTF8
WIN1253	<i>WIN1253</i> , UTF8
WIN1254	<i>WIN1254</i> , UTF8
WIN1255	<i>WIN1255</i> , UTF8
WIN1256	<i>WIN1256</i> , UTF8
WIN1257	<i>WIN1257</i> , UTF8
WIN1258	<i>WIN1258</i> , UTF8

Table 23.5. All Built-in Character Set Conversions

Conversion Name ^a	Source Encoding	Destination Encoding
big5_to_euc_tw	BIG5	EUC_TW
big5_to_mic	BIG5	MULE_INTERNAL
big5_to_utf8	BIG5	UTF8
euc_cn_to_mic	EUC_CN	MULE_INTERNAL
euc_cn_to_utf8	EUC_CN	UTF8
euc_jp_to_mic	EUC_JP	MULE_INTERNAL
euc_jp_to_sjis	EUC_JP	SJIS
euc_jp_to_utf8	EUC_JP	UTF8
euc_kr_to_mic	EUC_KR	MULE_INTERNAL
euc_kr_to_utf8	EUC_KR	UTF8
euc_tw_to_big5	EUC_TW	BIG5
euc_tw_to_mic	EUC_TW	MULE_INTERNAL
euc_tw_to_utf8	EUC_TW	UTF8
gb18030_to_utf8	GB18030	UTF8
gbk_to_utf8	GBK	UTF8
iso_8859_10_to_utf8	LATIN6	UTF8
iso_8859_13_to_utf8	LATIN7	UTF8
iso_8859_14_to_utf8	LATIN8	UTF8
iso_8859_15_to_utf8	LATIN9	UTF8
iso_8859_16_to_utf8	LATIN10	UTF8
iso_8859_1_to_mic	LATIN1	MULE_INTERNAL
iso_8859_1_to_utf8	LATIN1	UTF8
iso_8859_2_to_mic	LATIN2	MULE_INTERNAL
iso_8859_2_to_utf8	LATIN2	UTF8
iso_8859_2_to_windows_1250	LATIN2	WIN1250
iso_8859_3_to_mic	LATIN3	MULE_INTERNAL
iso_8859_3_to_utf8	LATIN3	UTF8
iso_8859_4_to_mic	LATIN4	MULE_INTERNAL
iso_8859_4_to_utf8	LATIN4	UTF8

Localization

Conversion Name ^a	Source Encoding	Destination Encoding
iso_8859_5_to_koi8_r	ISO_8859_5	KOI8R
iso_8859_5_to_mic	ISO_8859_5	MULE_INTERNAL
iso_8859_5_to_utf8	ISO_8859_5	UTF8
iso_8859_5_to_windows_1251	ISO_8859_5	WIN1251
iso_8859_5_to_windows_866	ISO_8859_5	WIN866
iso_8859_6_to_utf8	ISO_8859_6	UTF8
iso_8859_7_to_utf8	ISO_8859_7	UTF8
iso_8859_8_to_utf8	ISO_8859_8	UTF8
iso_8859_9_to_utf8	LATIN5	UTF8
johab_to_utf8	JOHAB	UTF8
koi8_r_to_iso_8859_5	KOI8R	ISO_8859_5
koi8_r_to_mic	KOI8R	MULE_INTERNAL
koi8_r_to_utf8	KOI8R	UTF8
koi8_r_to_windows_1251	KOI8R	WIN1251
koi8_r_to_windows_866	KOI8R	WIN866
koi8_u_to_utf8	KOI8U	UTF8
mic_to_big5	MULE_INTERNAL	BIG5
mic_to_euc_cn	MULE_INTERNAL	EUC_CN
mic_to_euc_jp	MULE_INTERNAL	EUC_JP
mic_to_euc_kr	MULE_INTERNAL	EUC_KR
mic_to_euc_tw	MULE_INTERNAL	EUC_TW
mic_to_iso_8859_1	MULE_INTERNAL	LATIN1
mic_to_iso_8859_2	MULE_INTERNAL	LATIN2
mic_to_iso_8859_3	MULE_INTERNAL	LATIN3
mic_to_iso_8859_4	MULE_INTERNAL	LATIN4
mic_to_iso_8859_5	MULE_INTERNAL	ISO_8859_5
mic_to_koi8_r	MULE_INTERNAL	KOI8R
mic_to_sjis	MULE_INTERNAL	SJIS
mic_to_windows_1250	MULE_INTERNAL	WIN1250
mic_to_windows_1251	MULE_INTERNAL	WIN1251
mic_to_windows_866	MULE_INTERNAL	WIN866
sjis_to_euc_jp	SJIS	EUC_JP
sjis_to_mic	SJIS	MULE_INTERNAL
sjis_to_utf8	SJIS	UTF8
windows_1258_to_utf8	WIN1258	UTF8
uhc_to_utf8	UHC	UTF8
utf8_to_big5	UTF8	BIG5
utf8_to_euc_cn	UTF8	EUC_CN
utf8_to_euc_jp	UTF8	EUC_JP
utf8_to_euc_kr	UTF8	EUC_KR

Localization

Conversion Name ^a	Source Encoding	Destination Encoding
utf8_to_euc_tw	UTF8	EUC_TW
utf8_to_gb18030	UTF8	GB18030
utf8_to_gbk	UTF8	GBK
utf8_to_iso_8859_1	UTF8	LATIN1
utf8_to_iso_8859_10	UTF8	LATIN6
utf8_to_iso_8859_13	UTF8	LATIN7
utf8_to_iso_8859_14	UTF8	LATIN8
utf8_to_iso_8859_15	UTF8	LATIN9
utf8_to_iso_8859_16	UTF8	LATIN10
utf8_to_iso_8859_2	UTF8	LATIN2
utf8_to_iso_8859_3	UTF8	LATIN3
utf8_to_iso_8859_4	UTF8	LATIN4
utf8_to_iso_8859_5	UTF8	ISO_8859_5
utf8_to_iso_8859_6	UTF8	ISO_8859_6
utf8_to_iso_8859_7	UTF8	ISO_8859_7
utf8_to_iso_8859_8	UTF8	ISO_8859_8
utf8_to_iso_8859_9	UTF8	LATIN5
utf8_to_johab	UTF8	JOHAB
utf8_to_koi8_r	UTF8	KOI8R
utf8_to_koi8_u	UTF8	KOI8U
utf8_to_sjis	UTF8	SJIS
utf8_to_windows_1258	UTF8	WIN1258
utf8_to_uhc	UTF8	UHC
utf8_to_windows_1250	UTF8	WIN1250
utf8_to_windows_1251	UTF8	WIN1251
utf8_to_windows_1252	UTF8	WIN1252
utf8_to_windows_1253	UTF8	WIN1253
utf8_to_windows_1254	UTF8	WIN1254
utf8_to_windows_1255	UTF8	WIN1255
utf8_to_windows_1256	UTF8	WIN1256
utf8_to_windows_1257	UTF8	WIN1257
utf8_to_windows_866	UTF8	WIN866
utf8_to_windows_874	UTF8	WIN874
windows_1250_to_iso_8859_2	WIN1250	LATIN2
windows_1250_to_mic	WIN1250	MULE_INTERNAL
windows_1250_to_utf8	WIN1250	UTF8
windows_1251_to_iso_8859_5	WIN1251	ISO_8859_5
windows_1251_to_koi8_r	WIN1251	KOI8R
windows_1251_to_mic	WIN1251	MULE_INTERNAL
windows_1251_to_utf8	WIN1251	UTF8

Conversion Name ^a	Source Encoding	Destination Encoding
windows_1251_to_windows_866	WIN1251	WIN866
windows_1252_to_utf8	WIN1252	UTF8
windows_1256_to_utf8	WIN1256	UTF8
windows_866_to_iso_8859_5	WIN866	ISO_8859_5
windows_866_to_koi8_r	WIN866	KOI8R
windows_866_to_mic	WIN866	MULE_INTERNAL
windows_866_to_utf8	WIN866	UTF8
windows_866_to_windows_1251	WIN866	WIN
windows_874_to_utf8	WIN874	UTF8
euc_jis_2004_to_utf8	EUC_JIS_2004	UTF8
utf8_to_euc_jis_2004	UTF8	EUC_JIS_2004
shift_jis_2004_to_utf8	SHIFT_JIS_2004	UTF8
utf8_to_shift_jis_2004	UTF8	SHIFT_JIS_2004
euc_jis_2004_to_shift_jis_2004	EUC_JIS_2004	SHIFT_JIS_2004
shift_jis_2004_to_euc_jis_2004	SHIFT_JIS_2004	EUC_JIS_2004

^a The conversion names follow a standard naming scheme: The official name of the source encoding with all non-alphanumeric characters replaced by underscores, followed by `_to_`, followed by the similarly processed destination encoding name. Therefore, these names sometimes deviate from the customary encoding names shown in Table 23.3.

23.3.5. Further Reading

These are good sources to start learning about various kinds of encoding systems.

CJKV Information Processing: Chinese, Japanese, Korean & Vietnamese Computing

Contains detailed explanations of EUC_JP, EUC_CN, EUC_KR, EUC_TW.

<https://www.unicode.org/>

The web site of the Unicode Consortium.

RFC 3629

UTF-8 (8-bit UCS/Unicode Transformation Format) is defined here.

Chapter 24. Routine Database Maintenance Tasks

Postgres Pro, like any database software, requires that certain tasks be performed regularly to achieve optimum performance. The tasks discussed here are *required*, but they are repetitive in nature and can easily be automated using standard tools such as cron scripts or Windows' Task Scheduler. It is the database administrator's responsibility to set up appropriate scripts, and to check that they execute successfully.

One obvious maintenance task is the creation of backup copies of the data on a regular schedule. Without a recent backup, you have no chance of recovery after a catastrophe (disk failure, fire, mistakenly dropping a critical table, etc.). The backup and recovery mechanisms available in Postgres Pro are discussed at length in [Chapter 25](#).

The other main category of maintenance task is periodic “vacuuming” of the database. This activity is discussed in [Section 24.1](#). Closely related to this is updating the statistics that will be used by the query planner, as discussed in [Section 24.1.3](#).

Another task that might need periodic attention is log file management. This is discussed in [Section 24.4](#).

[check_postgres](#) is available for monitoring database health and reporting unusual conditions. [check_postgres](#) integrates with Nagios and MRTG, but can be run standalone too.

Postgres Pro is low-maintenance compared to some other database management systems. Nonetheless, appropriate attention to these tasks will go far towards ensuring a pleasant and productive experience with the system.

24.1. Routine Vacuuming

Postgres Pro databases require periodic maintenance known as *vacuuming*. For many installations, it is sufficient to let vacuuming be performed by the *autovacuum daemon*, which is described in [Section 24.1.6](#). You might need to adjust the autovacuuming parameters described there to obtain best results for your situation. Some database administrators will want to supplement or replace the daemon's activities with manually-managed `VACUUM` commands, which typically are executed according to a schedule by cron or Task Scheduler scripts. To set up manually-managed vacuuming properly, it is essential to understand the issues discussed in the next few subsections. Administrators who rely on autovacuuming may still wish to skim this material to help them understand and adjust autovacuuming.

24.1.1. Vacuuming Basics

Postgres Pro's `VACUUM` command has to process each table on a regular basis for several reasons:

1. To recover or reuse disk space occupied by updated or deleted rows.
2. To update data statistics used by the Postgres Pro query planner.
3. To update the visibility map, which speeds up [index-only scans](#).
4. To shrink `pg_xact` and `pg_multixact`.

Each of these reasons dictates performing `VACUUM` operations of varying frequency and scope, as explained in the following subsections.

There are two variants of `VACUUM`: standard `VACUUM` and `VACUUM FULL`. `VACUUM FULL` can reclaim more disk space but runs much more slowly. Also, the standard form of `VACUUM` can run in parallel with production database operations. (Commands such as `SELECT`, `INSERT`, `UPDATE`, and `DELETE` will continue to function normally, though you will not be able to modify the definition of a table with commands such as `ALTER TABLE` while it is being vacuumed.) `VACUUM FULL` requires an `ACCESS EXCLUSIVE` lock on the table it is

working on, and therefore cannot be done in parallel with other use of the table. Generally, therefore, administrators should strive to use standard `VACUUM` and avoid `VACUUM FULL`.

`VACUUM` creates a substantial amount of I/O traffic, which can cause poor performance for other active sessions. There are configuration parameters that can be adjusted to reduce the performance impact of background vacuuming — see [Section 19.4.4](#).

24.1.2. Recovering Disk Space

In Postgres Pro, an `UPDATE` or `DELETE` of a row does not immediately remove the old version of the row. This approach is necessary to gain the benefits of multiversion concurrency control (MVCC, see [Chapter 13](#)): the row version must not be deleted while it is still potentially visible to other transactions. But eventually, an outdated or deleted row version is no longer of interest to any transaction. The space it occupies must then be reclaimed for reuse by new rows, to avoid unbounded growth of disk space requirements. This is done by running `VACUUM`.

The standard form of `VACUUM` removes dead row versions in tables and indexes and marks the space available for future reuse. However, it will not return the space to the operating system, except in the special case where one or more pages at the end of a table become entirely free and an exclusive table lock can be easily obtained. In contrast, `VACUUM FULL` actively compacts tables by writing a complete new version of the table file with no dead space. This minimizes the size of the table, but can take a long time. It also requires extra disk space for the new copy of the table, until the operation completes.

The usual goal of routine vacuuming is to do standard `VACUUMS` often enough to avoid needing `VACUUM FULL`. The autovacuum daemon attempts to work this way, and in fact will never issue `VACUUM FULL`. In this approach, the idea is not to keep tables at their minimum size, but to maintain steady-state usage of disk space: each table occupies space equivalent to its minimum size plus however much space gets used up between vacuum runs. Although `VACUUM FULL` can be used to shrink a table back to its minimum size and return the disk space to the operating system, there is not much point in this if the table will just grow again in the future. Thus, moderately-frequent standard `VACUUM` runs are a better approach than infrequent `VACUUM FULL` runs for maintaining heavily-updated tables.

Some administrators prefer to schedule vacuuming themselves, for example doing all the work at night when load is low. The difficulty with doing vacuuming according to a fixed schedule is that if a table has an unexpected spike in update activity, it may get bloated to the point that `VACUUM FULL` is really necessary to reclaim space. Using the autovacuum daemon alleviates this problem, since the daemon schedules vacuuming dynamically in response to update activity. It is unwise to disable the daemon completely unless you have an extremely predictable workload. One possible compromise is to set the daemon's parameters so that it will only react to unusually heavy update activity, thus keeping things from getting out of hand, while scheduled `VACUUMS` are expected to do the bulk of the work when the load is typical.

For those not using autovacuum, a typical approach is to schedule a database-wide `VACUUM` once a day during a low-usage period, supplemented by more frequent vacuuming of heavily-updated tables as necessary. (Some installations with extremely high update rates vacuum their busiest tables as often as once every few minutes.) If you have multiple databases in a cluster, don't forget to `VACUUM` each one; the program `vacuumdb` might be helpful.

Tip

Plain `VACUUM` may not be satisfactory when a table contains large numbers of dead row versions as a result of massive update or delete activity. If you have such a table and you need to reclaim the excess disk space it occupies, you will need to use `VACUUM FULL`, or alternatively `CLUSTER` or one of the table-rewriting variants of `ALTER TABLE`. These commands rewrite an entire new copy of the table and build new indexes for it. All these options require an `ACCESS EXCLUSIVE` lock. Note that they also temporarily use extra disk space approximately equal to the size of the table, since the old copies of the table and indexes can't be released until the new ones are complete.

Tip

If you have a table whose entire contents are deleted on a periodic basis, consider doing it with `TRUNCATE` rather than using `DELETE` followed by `VACUUM`. `TRUNCATE` removes the entire content of the table immediately, without requiring a subsequent `VACUUM` or `VACUUM FULL` to reclaim the now-unused disk space. The disadvantage is that strict MVCC semantics are violated.

24.1.1.3. Updating Planner Statistics

The Postgres Pro query planner relies on statistical information about the contents of tables in order to generate good plans for queries. These statistics are gathered by the `ANALYZE` command, which can be invoked by itself or as an optional step in `VACUUM`. It is important to have reasonably accurate statistics, otherwise poor choices of plans might degrade database performance.

The autovacuum daemon, if enabled, will automatically issue `ANALYZE` commands whenever the content of a table has changed sufficiently. However, administrators might prefer to rely on manually-scheduled `ANALYZE` operations, particularly if it is known that update activity on a table will not affect the statistics of “interesting” columns. The daemon schedules `ANALYZE` strictly as a function of the number of rows inserted or updated; it has no knowledge of whether that will lead to meaningful statistical changes.

Tuples changed in partitions and inheritance children do not trigger analyze on the parent table. If the parent table is empty or rarely changed, it may never be processed by autovacuum, and the statistics for the inheritance tree as a whole won't be collected. It is necessary to run `ANALYZE` on the parent table manually in order to keep the statistics up to date.

As with vacuuming for space recovery, frequent updates of statistics are more useful for heavily-updated tables than for seldom-updated ones. But even for a heavily-updated table, there might be no need for statistics updates if the statistical distribution of the data is not changing much. A simple rule of thumb is to think about how much the minimum and maximum values of the columns in the table change. For example, a `timestamp` column that contains the time of row update will have a constantly-increasing maximum value as rows are added and updated; such a column will probably need more frequent statistics updates than, say, a column containing URLs for pages accessed on a website. The URL column might receive changes just as often, but the statistical distribution of its values probably changes relatively slowly.

It is possible to run `ANALYZE` on specific tables and even just specific columns of a table, so the flexibility exists to update some statistics more frequently than others if your application requires it. In practice, however, it is usually best to just analyze the entire database, because it is a fast operation. `ANALYZE` uses a statistically random sampling of the rows of a table rather than reading every single row.

Tip

Although per-column tweaking of `ANALYZE` frequency might not be very productive, you might find it worthwhile to do per-column adjustment of the level of detail of the statistics collected by `ANALYZE`. Columns that are heavily used in `WHERE` clauses and have highly irregular data distributions might require a finer-grain data histogram than other columns. See `ALTER TABLE SET STATISTICS`, or change the database-wide default using the `default_statistics_target` configuration parameter.

Also, by default there is limited information available about the selectivity of functions. However, if you create a statistics object or an expression index that uses a function call, useful statistics will be gathered about the function, which can greatly improve query plans that use the expression index.

Tip

The autovacuum daemon does not issue `ANALYZE` commands for foreign tables, since it has no means of determining how often that might be useful. If your queries require statistics on foreign

tables for proper planning, it's a good idea to run manually-managed `ANALYZE` commands on those tables on a suitable schedule.

Tip

The autovacuum daemon does not issue `ANALYZE` commands for partitioned tables. Inheritance parents will only be analyzed if the parent itself is changed - changes to child tables do not trigger autoanalyze on the parent table. If your queries require statistics on parent tables for proper planning, it is necessary to periodically run a manual `ANALYZE` on those tables to keep the statistics up to date.

24.1.4. Updating the Visibility Map

Vacuum maintains a [visibility map](#) for each table to keep track of which pages contain only tuples that are known to be visible to all active transactions (and all future transactions, until the page is again modified). This has two purposes. First, vacuum itself can skip such pages on the next run, since there is nothing to clean up.

Second, it allows Postgres Pro to answer some queries using only the index, without reference to the underlying table. Since Postgres Pro indexes don't contain tuple visibility information, a normal index scan fetches the heap tuple for each matching index entry, to check whether it should be seen by the current transaction. An [index-only scan](#), on the other hand, checks the visibility map first. If it's known that all tuples on the page are visible, the heap fetch can be skipped. This is most useful on large data sets where the visibility map can prevent disk accesses. The visibility map is vastly smaller than the heap, so it can easily be cached even when the heap is very large.

24.1.5. Forced shrinking `pg_xact` and `pg_multixact`

Postgres Pro's [MVCC](#) transaction semantics depend on being able to compare transaction ID (XID) numbers: a row version with an insertion XID greater than the current transaction's XID is "in the future" and should not be visible to the current transaction.

Postgres Pro Enterprise 9.6 introduced 64-bit transaction IDs, which are not subject to wraparound and do not need modulo-2³² arithmetic to be compared. Each tuple header contains two XIDs, so extending them would lead to a significant overhead by occupying a lot of disk space. For that reason, when saved on disk on-tuple XIDs are 32-bit, but each page special area contains an offset, called *base XID*. For more information, see [Section 67.1](#).

The reason that periodic vacuuming solves the problem is that `VACUUM` will mark rows as *frozen*, indicating that they were inserted by a transaction that committed sufficiently far in the past that the effects of the inserting transaction are certain to be visible to all current and future transactions. Postgres Pro reserves a special XID, `FrozenTransactionId`, which is always considered older than every normal XID. Frozen row versions are treated as if the inserting XID were `FrozenTransactionId`, so that they will appear to be "in the past" to all normal transactions.

When a new XID can't fit into the existing page according to its base XID, this offset is shifted. Single page freezing takes place if needed. Both actions are performed "on the fly". Page-level wraparound can happen only when someone holds snapshot, which has more than 4 billions transactions. Even in highly loaded systems that may have up to 2³² transactions per day, it would take 2³¹ days before the first enforced "vacuum to prevent wraparound". The routine vacuuming is not enforced, and a database administrator can plan this operation during periods of low system load and minimal impact on database performance. Moreover, it can be done less frequently on systems with transaction rates similar to those mentioned above: several times a year instead of every few days.

Freezing data by `VACUUM` is needed to shrink `pg_xact` and `pg_multixact`. For historical reasons, the wording "autovacuum to prevent wraparound" is preserved for forced autovacuum shrinking `pg_xact` and `pg_multixact`.

Note

In PostgreSQL versions before 9.4, freezing was implemented by actually replacing a row's insertion `XID` with `FrozenTransactionId`, which was visible in the row's `xmin` system column. Newer versions just set a flag bit, preserving the row's original `xmin` for possible forensic use. However, rows with `xmin` equal to `FrozenTransactionId` (2) may still be found in databases `pg_upgrade'd` from pre-9.4 versions.

Also, system catalogs may contain rows with `xmin` equal to `BootstrapTransactionId` (1), indicating that they were inserted during the first phase of `initdb`. Like `FrozenTransactionId`, this special `XID` is treated as older than every normal `XID`.

`vacuum_freeze_min_age` controls how old an `XID` value has to be before rows bearing that `XID` will be frozen. Increasing this setting may avoid unnecessary work if the rows that would otherwise be frozen will soon be modified again, but decreasing this setting increases the number of transactions that can elapse before the table must be vacuumed again.

`VACUUM` uses the [visibility map](#) to determine which pages of a table must be scanned. Normally, it will skip pages that don't have any dead row versions even if those pages might still have row versions with old `XID` values. Therefore, normal `VACUUMS` won't always freeze every old row version in the table. When that happens, `VACUUM` will eventually need to perform an *aggressive vacuum*, which will freeze all eligible unfrozen `XID` and `MXID` values, including those from all-visible but not all-frozen pages. In practice most tables require periodic aggressive vacuuming. `vacuum_freeze_table_age` controls when `VACUUM` does that: all-visible but not all-frozen pages are scanned if the number of transactions that have passed since the last such scan is greater than `vacuum_freeze_table_age` minus `vacuum_freeze_min_age`. Setting `vacuum_freeze_table_age` to 0 forces `VACUUM` to always use its aggressive strategy.

Autovacuum is invoked on any table that might contain unfrozen rows with `XIDs` older than the age specified by the configuration parameter `autovacuum_freeze_max_age`. (This will happen even if autovacuum is disabled.)

This implies that if a table is not otherwise vacuumed, autovacuum will be invoked on it approximately once every `autovacuum_freeze_max_age` minus `vacuum_freeze_min_age` transactions. For tables that are regularly vacuumed for space reclamation purposes, this is of little importance. However, for static tables (including tables that receive inserts, but no updates or deletes), there is no need to vacuum for space reclamation, so it can be useful to try to maximize the interval between forced autovacuaums on very large static tables. Obviously one can do this either by increasing `autovacuum_freeze_max_age` or decreasing `vacuum_freeze_min_age`.

The effective maximum for `vacuum_freeze_table_age` is $0.95 * \text{autovacuum_freeze_max_age}$; a setting higher than that will be capped to the maximum. A value higher than `autovacuum_freeze_max_age` wouldn't make sense because an autovacuum to shrink `pg_xact` and `pg_multixact` would be triggered at that point anyway, and the 0.95 multiplier leaves some breathing room to run a manual `VACUUM` before that happens. As a rule of thumb, `vacuum_freeze_table_age` should be set to a value somewhat below `autovacuum_freeze_max_age`, leaving enough gap so that a regularly scheduled `VACUUM` or an autovacuum triggered by normal delete and update activity is run in that window. Setting it too close could lead to autovacuaums to shrink `pg_xact` and `pg_multixact`, even though the table was recently vacuumed to reclaim space, whereas lower values lead to more frequent aggressive vacuuming.

The sole disadvantage of increasing `autovacuum_freeze_max_age` (and `vacuum_freeze_table_age` along with it) is that the `pg_xact` and `pg_commit_ts` subdirectories of the database cluster will take more space, because it must store the commit status and (if `track_commit_timestamp` is enabled) timestamp of all transactions back to the `autovacuum_freeze_max_age` horizon. The commit status uses two bits per transaction, so if `autovacuum_freeze_max_age` is set to its maximum allowed value of two billion, `pg_xact` can be expected to grow to about half a gigabyte and `pg_commit_ts` to about 20GB. If this is trivial compared to your total database size, setting `autovacuum_freeze_max_age` to its maximum al-

lowed value is recommended. Otherwise, set it depending on what you are willing to allow for `pg_xact` and `pg_commit_ts` storage. (The default, 200 million transactions, translates to about 50MB of `pg_xact` storage and about 2GB of `pg_commit_ts` storage.)

One disadvantage of decreasing `vacuum_freeze_min_age` is that it might cause `VACUUM` to do useless work: freezing a row version is a waste of time if the row is modified soon thereafter (causing it to acquire a new XID). So the setting should be large enough that rows are not frozen until they are unlikely to change any more.

To track the age of the oldest unfrozen XIDs in a database, `VACUUM` stores XID statistics in the system tables `pg_class` and `pg_database`. In particular, the `relfrozenxid` column of a table's `pg_class` row contains the oldest remaining unfrozen XID at the end of the most recent `VACUUM` that successfully advanced `relfrozenxid` (typically the most recent aggressive `VACUUM`). Similarly, the `datfrozenxid` column of a database's `pg_database` row is a lower bound on the unfrozen XIDs appearing in that database — it is just the minimum of the per-table `relfrozenxid` values within the database. A convenient way to examine this information is to execute queries such as:

```
SELECT c.oid::regclass as table_name,
       greatest(age(c.relfrozenxid), age(t.relfrozenxid)) as age
FROM pg_class c
LEFT JOIN pg_class t ON c.reltoastrelid = t.oid
WHERE c.relkind IN ('r', 'm');

SELECT datname, age(datfrozenxid) FROM pg_database;
```

The `age` column measures the number of transactions from the cutoff XID to the current transaction's XID.

Tip

When the `VACUUM` command's `VERBOSE` parameter is specified, `VACUUM` prints various statistics about the table. This includes information about how `relfrozenxid` and `relminmxid` advanced, and the number of newly frozen pages. The same details appear in the server log when `autovacuum` logging (controlled by [log_autovacuum_min_duration](#)) reports on a `VACUUM` operation executed by `autovacuum`.

`VACUUM` normally only scans pages that have been modified since the last vacuum, but `relfrozenxid` can only be advanced when every page of the table that might contain unfrozen XIDs is scanned. This happens when `relfrozenxid` is more than `vacuum_freeze_table_age` transactions old, when `VACUUM`'s `FREEZE` option is used, or when all pages that are not already all-frozen happen to require vacuuming to remove dead row versions. When `VACUUM` scans every page in the table that is not already all-frozen, it should set `age(relfrozenxid)` to a value just a little more than the `vacuum_freeze_min_age` setting that was used (more by the number of transactions started since the `VACUUM` started). `VACUUM` will set `relfrozenxid` to the oldest XID that remains in the table, so it's possible that the final value will be much more recent than strictly required. If no `relfrozenxid`-advancing `VACUUM` is issued on the table until `autovacuum_freeze_max_age` is reached, an `autovacuum` will soon be forced for the table.

24.1.5.1. Shrinking `pg_multixact`

Multixact IDs are used to support row locking by multiple transactions. Since there is only limited space in a tuple header to store lock information, that information is encoded as a “multiple transaction ID”, or *multixact ID* for short, whenever there is more than one transaction concurrently locking a row. Information about which transaction IDs are included in any particular multixact ID is stored separately in the `pg_multixact` subdirectory, and only the multixact ID appears in the `xmax` field in the tuple header. Like transaction IDs, multixact IDs are implemented on disk page as a 64-bit counter with an offset, and corresponding storage, which requires careful aging management, and storage cleanup. There is a separate storage area which holds the list of members in each multixact, which uses a 64-bit counter.

Whenever `VACUUM` scans any part of a table, it will replace any multixact ID it encounters which is older than `vacuum_multixact_freeze_min_age` by a different value, which can be the zero value, a single transaction ID, or a newer multixact ID. For each table, `pg_class.relminmxid` stores the oldest possible multixact ID still appearing in any tuple of that table. If this value is older than `vacuum_multixact_freeze_table_age`, an aggressive vacuum is forced. As discussed in the previous section, an aggressive vacuum means that only those pages which are known to be all-frozen will be skipped. `mxid_age()` can be used on `pg_class.relminmxid` to find its age.

Aggressive `VACUUMS`, regardless of what causes them, are *guaranteed* to be able to advance the table's `relminmxid`. Eventually, as all tables in all databases are scanned and their oldest multixact values are advanced, on-disk storage for older multixacts can be removed.

An aggressive vacuum scan will occur for any table whose multixact-age is greater than `autovacuum_multixact_freeze_max_age`.

24.1.6. The Autovacuum Daemon

Postgres Pro has an optional but highly recommended feature called *autovacuum*, whose purpose is to automate the execution of `VACUUM` and `ANALYZE` commands. When enabled, autovacuum checks for tables that have had a large number of inserted, updated or deleted tuples. These checks use the statistics collection facility; therefore, autovacuum cannot be used unless `track_counts` is set to `true`. In the default configuration, autovacuuming is enabled and the related configuration parameters are appropriately set.

The “autovacuum daemon” actually consists of multiple processes. There is a persistent daemon process, called the *autovacuum launcher*, which is in charge of starting *autovacuum worker* processes for all databases. The launcher will distribute the work across time, attempting to start one worker within each database every `autovacuum_naptime` seconds. (Therefore, if the installation has *N* databases, a new worker will be launched every `autovacuum_naptime/N` seconds.) A maximum of `autovacuum_max_workers` worker processes are allowed to run at the same time. If there are more than `autovacuum_max_workers` databases to be processed, the next database will be processed as soon as the first worker finishes. Each worker process will check each table within its database and execute `VACUUM` and/or `ANALYZE` as needed. `log_autovacuum_min_duration` can be set to monitor autovacuum workers' activity.

If an autovacuum worker process comes across a table with `parallel_autovacuum_workers` ≥ 0 , this worker tells the autovacuum launcher of the need to vacuum indexes of this table in a parallel mode. To this end, the launcher starts as many additional autovacuum worker processes as defined by `parallel_autovacuum_workers`. Parallel autovacuum is currently in an experimental phase and is intended for testing purposes only.

If several large tables all become eligible for vacuuming in a short amount of time, all autovacuum workers might become occupied with vacuuming those tables for a long period. This would result in other tables and databases not being vacuumed until a worker becomes available. There is no limit on how many workers might be in a single database, but workers do try to avoid repeating work that has already been done by other workers. Note that the number of running workers does not count towards `max_connections` or `superuser_reserved_connections` limits.

Tables whose `relfrozenxid` value is more than `autovacuum_freeze_max_age` transactions old are always vacuumed (this also applies to those tables whose freeze max age has been modified via storage parameters; see below). Otherwise, if the number of tuples obsoleted since the last `VACUUM` exceeds the “vacuum threshold”, the table is vacuumed. The vacuum threshold is defined as:

`vacuum threshold = vacuum base threshold + vacuum scale factor * number of tuples`

where the vacuum base threshold is `autovacuum_vacuum_threshold`, the vacuum scale factor is `autovacuum_vacuum_scale_factor`, and the number of tuples is `pg_class.reltuples`.

The table is also vacuumed if the number of tuples inserted since the last vacuum has exceeded the defined insert threshold, which is defined as:

`vacuum insert threshold = vacuum base insert threshold + vacuum insert scale factor * number of tuples`

where the vacuum insert base threshold is `autovacuum_vacuum_insert_threshold`, and vacuum insert scale factor is `autovacuum_vacuum_insert_scale_factor`. Such vacuums may allow portions of the table to be marked as *all visible* and also allow tuples to be frozen, which can reduce the work required in subsequent vacuums. For tables which receive `INSERT` operations but no or almost no `UPDATE/DELETE` operations, it may be beneficial to lower the table's `autovacuum_freeze_min_age` as this may allow tuples to be frozen by earlier vacuums. The number of obsolete tuples and the number of inserted tuples are obtained from the cumulative statistics system; it is an eventually-consistent count updated by each `UPDATE`, `DELETE` and `INSERT` operation. If the `relfrozenxid` value of the table is more than `vacuum_freeze_table_age` transactions old, an aggressive vacuum is performed to freeze old tuples and advance `relfrozenxid`; otherwise, only pages that have been modified since the last vacuum are scanned.

For analyze, a similar condition is used: the threshold, defined as:

```
analyze threshold = analyze base threshold + analyze scale factor * number of tuples
```

is compared to the total number of tuples inserted, updated, or deleted since the last `ANALYZE`.

Partitioned tables do not directly store tuples and consequently are not processed by autovacuum. (Autovacuum does process table partitions just like other tables.) Unfortunately, this means that autovacuum does not run `ANALYZE` on partitioned tables, and this can cause suboptimal plans for queries that reference partitioned table statistics. You can work around this problem by manually running `ANALYZE` on partitioned tables when they are first populated, and again whenever the distribution of data in their partitions changes significantly.

Temporary tables cannot be accessed by autovacuum. Therefore, appropriate vacuum and analyze operations should be performed via session SQL commands.

The default thresholds and scale factors are taken from `postgresql.conf`, but it is possible to override them (and many other autovacuum control parameters) on a per-table basis; see [Storage Parameters](#) for more information. If a setting has been changed via a table's storage parameters, that value is used when processing that table; otherwise the global settings are used. See [Section 19.10](#) for more details on the global settings.

When multiple workers are running, the autovacuum cost delay parameters (see [Section 19.4.4](#)) are “balanced” among all the running workers, so that the total I/O impact on the system is the same regardless of the number of workers actually running. However, any workers processing tables whose per-table `autovacuum_vacuum_cost_delay` or `autovacuum_vacuum_cost_limit` storage parameters have been set are not considered in the balancing algorithm.

Autovacuum workers generally don't block other commands. If a process attempts to acquire a lock that conflicts with the `SHARE UPDATE EXCLUSIVE` lock held by autovacuum, lock acquisition will interrupt the autovacuum. For conflicting lock modes, see [Table 13.2](#). However, if the autovacuum is running to prevent transaction ID wraparound (i.e., the autovacuum query name in the `pg_stat_activity` view ends with `(to prevent wraparound)`), the autovacuum is not automatically interrupted.

Warning

Regularly running commands that acquire locks conflicting with a `SHARE UPDATE EXCLUSIVE` lock (e.g., `ANALYZE`) can effectively prevent autovacums from ever completing.

24.1.7. Analyzing and Vacuuming in Distributed System

24.1.7.1. The Autovacuum Daemon

For many installations, it is sufficient to let vacuuming be performed by the autovacuum daemon. As in PostgreSQL installations, the autovacuum daemon will issue `ANALYZE` commands automatically whenever the content of a table has changed sufficiently. The `ANALYZE` command transfers statistics for foreign partitions from remote nodes.

Rebalance process can move or copy data between cluster nodes. After this operation, all transferred objects are analyzed automatically. As usual, local statistics are gathered, and remote statistics are fetched from foreign servers.

24.1.7.2. Analyzing the Whole Database

When running for the entire database, the `ANALYZE` command calculates statistics for local partitions of sharded tables on the current node and for global tables. It also transfers statistics for foreign partitions available on remote nodes. However, note that remote nodes may lack the statistics.

Still, running `ANALYZE` on each node may not give you full cluster-wide statistics, as they may have not been calculated on some nodes. Also, global tables are analyzed multiple times (for each node).

Instead, use the `shardman.global_analyze()` function to perform cluster-wide analysis. First, this function executes `ANALYZE` on each node for all local partitions of sharded tables, then sends these statistics to other nodes. Next, it selects one node per a global table and runs `ANALYZE` of this table on the selected node. Gathered statistics are broadcast to all other nodes in the cluster.

24.1.7.3. Analyzing Global Tables

When `ANALYZE` is run on a global table, only statistics for the corresponding local table on the current node are updated.

24.1.7.4. Analyzing Sharded Tables

When `ANALYZE` is run on a sharded table, statistics on its local partitions are updated, statistics for foreign partitions are transferred from remote nodes, if remote nodes have it. If statistics are missing on remote nodes, they are not calculated. If you run `ANALYZE` for the local table, only its statistics are updated.

When `ANALYZE` is run on a foreign table directly and a remote node does not have any statistics for the corresponding local table, this local table is analyzed remotely. Then statistics are transferred from the remote node.

24.1.7.5. Vacuuming

The `VACUUM` command running for the whole database performs vacuuming on the current node.

You can broadcast this command to perform cluster-wide vacuuming on all nodes. Set the `shardman.broadcast_ddl` configuration parameter to `on`, run `VACUUM`, and set `shardman.broadcast_ddl` to `off`.

When `VACUUM` is run on a sharded or global table, the statement is broadcast automatically. On all nodes, the command vacuums all local partitions of the sharded table and all local tables of the global table.

24.2. Background Freezer

When obsolete tuple versions rapidly accumulate in database tables and `autovacuum` cannot keep up with the load, you can use the Postgres Pro feature called *background freezer*. It is a lightweight alternative to `autovacuum` that works only with modified and non-locked data pages of permanent tables stored in RAM. Background freezer rapidly vacuums table pages in RAM, which allows adding new tuples. This reduces table growth and improves `autovacuum` efficiency.

Note

You do not need to use background freezer if `autovacuum` can cope with the load.

24.2.1. Overview

Postgres Pro launches a dedicated background freezer worker. All backends communicate with the worker using a UDP socket to send information about recently modified pages and completed transactions to this worker. This information is accumulated in special lists individual for each table. When the number

of modified pages in the list exceeds the threshold set with the `bgfreezer_pages_per_rel_to_clean` configuration parameter, this list is queued for vacuuming and freezing. The worker starts processing this queue after being notified that the transaction has completed. While the list of pages is being vacuumed, the worker acquires a read lock on the corresponding table. This prevents parallel operations, such as dropping or truncating the table, during the background freezer activity. For each page of the list, the background freezer worker vacuums and freezes the tuples that meet the relevant conditions. For more information, see [Section 24.1.6](#). If after vacuuming and freezing all page tuples are visible and frozen and the relevant [visibility map](#) page is in memory, background freezer sets the visibility and freezing bits on the visibility map page. After deleted tuples are vacuummed on the data page, the [free space map \(FSM\)](#) data is also adjusted. Without such an adjustment, vacuummed pages will not be reused.

The background freezer mechanism only applies to permanent table pages, leaving temporary tables unaffected.

Background freezer helps with issues that may occur when a large number of temporary tables is frequently created and dropped. In such cases, autovacuum may not be fast enough to vacuum system tables causing them to bloat. This primarily applies to `pg_type`, `pg_attribute`, `pg_class`, `pg_depend`, and `pg_index`. Background freezer vacuums deleted tuples in memory, thereby considerably decreasing system table bloat.

To enable background freezer, set the `enable_background_freezer` configuration parameter to `on`.

24.3. Routine Reindexing

In some situations it is worthwhile to rebuild indexes periodically with the `REINDEX` command or a series of individual rebuilding steps.

B-tree index pages that have become completely empty are reclaimed for re-use. However, there is still a possibility of inefficient use of space: if all but a few index keys on a page have been deleted, the page remains allocated. Therefore, a usage pattern in which most, but not all, keys in each range are eventually deleted will see poor use of space. For such usage patterns, periodic reindexing is recommended.

The potential for bloat in non-B-tree indexes has not been well researched. It is a good idea to periodically monitor the index's physical size when using any non-B-tree index type.

Also, for B-tree indexes, a freshly-constructed index is slightly faster to access than one that has been updated many times because logically adjacent pages are usually also physically adjacent in a newly built index. (This consideration does not apply to non-B-tree indexes.) It might be worthwhile to reindex periodically just to improve access speed.

`REINDEX` can be used safely and easily in all cases. This command requires an `ACCESS EXCLUSIVE` lock by default, hence it is often preferable to execute it with its `CONCURRENTLY` option, which requires only a `SHARE UPDATE EXCLUSIVE` lock.

24.4. Log File Maintenance

It is a good idea to save the database server's log output somewhere, rather than just discarding it via `/dev/null`. The log output is invaluable when diagnosing problems.

Note

The server log can contain sensitive information and needs to be protected, no matter how or where it is stored, or the destination to which it is routed. For example, some DDL statements might contain plaintext passwords or other authentication details. Logged statements at the `ERROR` level might show the SQL source code for applications and might also contain some parts of data rows. Recording data, events and related information is the intended function of this facility, so this is not a leakage or a bug. Please ensure the server logs are visible only to appropriately authorized people.

Log output tends to be voluminous (especially at higher debug levels) so you won't want to save it indefinitely. You need to *rotate* the log files so that new log files are started and old ones removed after a reasonable period of time.

If you simply direct the stderr of `postgres` into a file, you will have log output, but the only way to truncate the log file is to stop and restart the server. This might be acceptable if you are using Postgres Pro in a development environment, but few production servers would find this behavior acceptable.

A better approach is to send the server's stderr output to some type of log rotation program. There is a built-in log rotation facility, which you can use by setting the configuration parameter `logging_collector` to `true` in `postgresql.conf`. The control parameters for this program are described in [Section 19.8.1](#). You can also use this approach to capture the log data in machine readable CSV (comma-separated values) format.

Alternatively, you might prefer to use an external log rotation program if you have one that you are already using with other server software. For example, the `rotatelogs` tool included in the Apache distribution can be used with Postgres Pro. One way to do this is to pipe the server's stderr output to the desired program. If you start the server with `pg_ctl`, then stderr is already redirected to stdout, so you just need a pipe command, for example:

```
pg_ctl start | rotatelogs /var/log/pgsql_log 86400
```

You can combine these approaches by setting up `logrotate` to collect log files produced by Postgres Pro built-in logging collector. In this case, the logging collector defines the names and location of the log files, while `logrotate` periodically archives these files. When initiating log rotation, `logrotate` must ensure that the application sends further output to the new file. This is commonly done with a `postrotate` script that sends a `SIGHUP` signal to the application, which then reopens the log file. In Postgres Pro, you can run `pg_ctl` with the `logrotate` option instead. When the server receives this command, the server either switches to a new log file or reopens the existing file, depending on the logging configuration (see [Section 19.8.1](#)).

Note

When using static log file names, the server might fail to reopen the log file if the max open file limit is reached or a file table overflow occurs. In this case, log messages are sent to the old log file until a successful log rotation. If `logrotate` is configured to compress the log file and delete it, the server may lose the messages logged in this time frame. To avoid this issue, you can configure the logging collector to dynamically assign log file names and use a `prerotate` script to ignore open log files.

Another production-grade approach to managing log output is to send it to syslog and let syslog deal with file rotation. To do this, set the configuration parameter `log_destination` to `syslog` (to log to syslog only) in `postgresql.conf`. Then you can send a `SIGHUP` signal to the syslog daemon whenever you want to force it to start writing a new log file. If you want to automate log rotation, the `logrotate` program can be configured to work with log files from syslog.

On many systems, however, syslog is not very reliable, particularly with large log messages; it might truncate or drop messages just when you need them the most. Also, on Linux, syslog will flush each message to disk, yielding poor performance. (You can use a `-` at the start of the file name in the syslog configuration file to disable syncing.)

Note that all the solutions described above take care of starting new log files at configurable intervals, but they do not handle deletion of old, no-longer-useful log files. You will probably want to set up a batch job to periodically delete old log files. Another possibility is to configure the rotation program so that old log files are overwritten cyclically.

Chapter 25. Backup and Restore

As with everything that contains valuable data, Postgres Pro databases should be backed up regularly. While the procedure is essentially simple, it is important to have a clear understanding of the underlying techniques and assumptions.

25.1. Continuous Archiving and Point-in-Time Recovery (PITR)

At all times, Postgres Pro maintains a *write ahead log* (WAL) in the `pg_wal/` subdirectory of the cluster's data directory. The log records every change made to the database's data files. This log exists primarily for crash-safety purposes: if the system crashes, the database can be restored to consistency by “replaying” the log entries made since the last checkpoint. However, the existence of the log makes it possible to use a third strategy for backing up databases: we can combine a file-system-level backup with backup of the WAL files. If recovery is needed, we restore the file system backup and then replay from the backed-up WAL files to bring the system to a current state. This approach is more complex to administer than either of the previous approaches, but it has some significant benefits:

- We do not need a perfectly consistent file system backup as the starting point. Any internal inconsistency in the backup will be corrected by log replay (this is not significantly different from what happens during crash recovery). So we do not need a file system snapshot capability, just tar or a similar archiving tool.
- Since we can combine an indefinitely long sequence of WAL files for replay, continuous backup can be achieved simply by continuing to archive the WAL files. This is particularly valuable for large databases, where it might not be convenient to take a full backup frequently.
- It is not necessary to replay the WAL entries all the way to the end. We could stop the replay at any point and have a consistent snapshot of the database as it was at that time. Thus, this technique supports *point-in-time recovery*: it is possible to restore the database to its state at any time since your base backup was taken.
- If we continuously feed the series of WAL files to another machine that has been loaded with the same base backup file, we have a *warm standby* system: at any point we can bring up the second machine and it will have a nearly-current copy of the database.

Note

`pg_dump` and `pg_dumpall` do not produce file-system-level backups and cannot be used as part of a continuous-archiving solution. Such dumps are *logical* and do not contain enough information to be used by WAL replay.

As with the plain file-system-backup technique, this method can only support restoration of an entire database cluster, not a subset. Also, it requires a lot of archival storage: the base backup might be bulky, and a busy system will generate many megabytes of WAL traffic that have to be archived. Still, it is the preferred backup technique in many situations where high reliability is needed.

To recover successfully using continuous archiving (also called “online backup” by many database vendors), you need a continuous sequence of archived WAL files that extends back at least as far as the start time of your backup. So to get started, you should set up and test your procedure for archiving WAL files *before* you take your first base backup. Accordingly, here we discuss some aspects of archiving WAL files. See [pg_probackup](#) for further reading.

25.1.1. Setting Up WAL Archiving

In an abstract sense, a running Postgres Pro system produces an indefinitely long sequence of WAL records. The system physically divides this sequence into WAL *segment files*, which are normally 16MB apiece (although the segment size can be altered during `initdb`). The segment files are given numeric

names that reflect their position in the abstract WAL sequence. When not using WAL archiving, the system normally creates just a few segment files and then “recycles” them by renaming no-longer-needed segment files to higher segment numbers. It's assumed that segment files whose contents precede the last checkpoint are no longer of interest and can be recycled.

When archiving WAL data, we need to capture the contents of each segment file once it is filled, and save that data somewhere before the segment file is recycled for reuse. Depending on the application and the available hardware, there could be many different ways of “saving the data somewhere”: we could copy the segment files to an NFS-mounted directory on another machine, write them onto a tape drive (ensuring that you have a way of identifying the original name of each file), or batch them together and burn them onto CDs, or something else entirely. To provide the database administrator with flexibility, Postgres Pro tries not to make any assumptions about how the archiving will be done. Instead, Postgres Pro lets the administrator specify a shell command or an archive library to be executed to copy a completed segment file to wherever it needs to go. This could be as simple as a shell command that uses `cp`, or it could invoke a complex C function — it's all up to you.

To enable WAL archiving, set the `wal_level` configuration parameter to `replica` or higher, `archive_mode` to `on`, specify the shell command to use in the `archive_command` configuration parameter or specify the library to use in the `archive_library` configuration parameter. In practice these settings will always be placed in the `postgresql.conf` file.

In `archive_command`, `%p` is replaced by the path name of the file to archive, while `%f` is replaced by only the file name. (The path name is relative to the current working directory, i.e., the cluster's data directory.) Use `%%` if you need to embed an actual `%` character in the command. The simplest useful command is something like:

```
archive_command = 'test ! -f /mnt/server/archivedir/%f && cp %p /mnt/server/archivedir/%f' # Unix
archive_command = 'copy "%p" "C:\\server\\archivedir\\%f"' # Windows
```

which will copy archivable WAL segments to the directory `/mnt/server/archivedir`. (This is an example, not a recommendation, and might not work on all platforms.) After the `%p` and `%f` parameters have been replaced, the actual command executed might look like this:

```
test ! -f /mnt/server/archivedir/000000010000000A9000000065 && cp
pg_wal/000000010000000A9000000065 /mnt/server/archivedir/000000010000000A9000000065
```

A similar command will be generated for each new file to be archived.

The archive command will be executed under the ownership of the same user that the Postgres Pro server is running as. Since the series of WAL files being archived contains effectively everything in your database, you will want to be sure that the archived data is protected from prying eyes; for example, archive into a directory that does not have group or world read access.

It is important that the archive command return zero exit status if and only if it succeeds. Upon getting a zero result, Postgres Pro will assume that the file has been successfully archived, and will remove or recycle it. However, a nonzero status tells Postgres Pro that the file was not archived; it will try again periodically until it succeeds.

Another way to archive is to use a custom archive module as the `archive_library`. Since such modules are written in C, creating your own may require considerably more effort than writing a shell command. However, archive modules can be more performant than archiving via shell, and they will have access to many useful server resources. For more information about archive modules, see [Chapter 52](#).

When the archive command is terminated by a signal (other than SIGTERM that is used as part of a server shutdown) or an error by the shell with an exit status greater than 125 (such as command not found), or if the archive function emits an `ERROR` or `FATAL`, the archiver process aborts and gets restarted by the postmaster. In such cases, the failure is not reported in `pg_stat_archiver`.

Archive commands and libraries should generally be designed to refuse to overwrite any pre-existing archive file. This is an important safety feature to preserve the integrity of your archive in case of ad-

ministrator error (such as sending the output of two different servers to the same archive directory). It is advisable to test your proposed archive library to ensure that it does not overwrite an existing file.

In rare cases, Postgres Pro may attempt to re-archive a WAL file that was previously archived. For example, if the system crashes before the server makes a durable record of archival success, the server will attempt to archive the file again after restarting (provided archiving is still enabled). When an archive command or library encounters a pre-existing file, it should return a zero status or `true`, respectively, if the WAL file has identical contents to the pre-existing archive and the pre-existing archive is fully persisted to storage. If a pre-existing file contains different contents than the WAL file being archived, the archive command or library *must* return a nonzero status or `false`, respectively.

The example command above for Unix avoids overwriting a pre-existing archive by including a separate `test` step. On some Unix platforms, `cp` has switches such as `-i` that can be used to do the same thing less verbosely, but you should not rely on these without verifying that the right exit status is returned. (In particular, GNU `cp` will return status zero when `-i` is used and the target file already exists, which is *not* the desired behavior.)

While designing your archiving setup, consider what will happen if the archive command or library fails repeatedly because some aspect requires operator intervention or the archive runs out of space. For example, this could occur if you write to tape without an autochanger; when the tape fills, nothing further can be archived until the tape is swapped. You should ensure that any error condition or request to a human operator is reported appropriately so that the situation can be resolved reasonably quickly. The `pg_wal/` directory will continue to fill with WAL segment files until the situation is resolved. (If the file system containing `pg_wal/` fills up, Postgres Pro will do a PANIC shutdown. No committed transactions will be lost, but the database will remain offline until you free some space.)

The speed of the archive command or library is unimportant as long as it can keep up with the average rate at which your server generates WAL data. Normal operation continues even if the archiving process falls a little behind. If archiving falls significantly behind, this will increase the amount of data that would be lost in the event of a disaster. It will also mean that the `pg_wal/` directory will contain large numbers of not-yet-archived segment files, which could eventually exceed available disk space. You are advised to monitor the archiving process to ensure that it is working as you intend.

In writing your archive command or library, you should assume that the file names to be archived can be up to 64 characters long and can contain any combination of ASCII letters, digits, and dots. It is not necessary to preserve the original relative path (`%p`) but it is necessary to preserve the file name (`%f`).

Note that although WAL archiving will allow you to restore any modifications made to the data in your Postgres Pro database, it will not restore changes made to configuration files (that is, `postgresql.conf`, `pg_hba.conf` and `pg_ident.conf`), since those are edited manually rather than through SQL operations. You might wish to keep the configuration files in a location that will be backed up by your regular file system backup procedures. See [Section 19.2](#) for how to relocate the configuration files.

The archive command or function is only invoked on completed WAL segments. Hence, if your server generates only little WAL traffic (or has slack periods where it does so), there could be a long delay between the completion of a transaction and its safe recording in archive storage. To put a limit on how old unarchived data can be, you can set `archive_timeout` to force the server to switch to a new WAL segment file at least that often. Note that archived files that are archived early due to a forced switch are still the same length as completely full files. It is therefore unwise to set a very short `archive_timeout` — it will bloat your archive storage. `archive_timeout` settings of a minute or so are usually reasonable.

Also, you can force a segment switch manually with `pg_switch_wal` if you want to ensure that a just-finished transaction is archived as soon as possible. Other utility functions related to WAL management are listed in [Table 9.95](#).

When `wal_level` is `minimal` some SQL commands are optimized to avoid WAL logging. If archiving or streaming replication were turned on during execution of one of these statements, WAL would not contain enough information for archive recovery. (Crash recovery is unaffected.) For this reason, `wal_level` can only be changed at server start. However, `archive_command` and `archive_library` can be changed

with a configuration file reload. If you are archiving via shell and wish to temporarily stop archiving, one way to do it is to set `archive_command` to the empty string (`' '`). This will cause WAL files to accumulate in `pg_wal/` until a working `archive_command` is re-established.

25.1.2. Timelines

The ability to restore the database to a previous point in time creates some complexities that are akin to science-fiction stories about time travel and parallel universes. For example, in the original history of the database, suppose you dropped a critical table at 5:15PM on Tuesday evening, but didn't realize your mistake until Wednesday noon. Unfazed, you get out your backup, restore to the point-in-time 5:14PM Tuesday evening, and are up and running. In *this* history of the database universe, you never dropped the table. But suppose you later realize this wasn't such a great idea, and would like to return to sometime Wednesday morning in the original history. You won't be able to if, while your database was up-and-running, it overwrote some of the WAL segment files that led up to the time you now wish you could get back to. Thus, to avoid this, you need to distinguish the series of WAL records generated after you've done a point-in-time recovery from those that were generated in the original database history.

To deal with this problem, Postgres Pro has a notion of *timelines*. Whenever an archive recovery completes, a new timeline is created to identify the series of WAL records generated after that recovery. The timeline ID number is part of WAL segment file names so a new timeline does not overwrite the WAL data generated by previous timelines. For example, in the WAL file name `0000000100001234000055CD`, the leading `00000001` is the timeline ID in hexadecimal. (Note that in other contexts, such as server log messages, timeline IDs are usually printed in decimal.)

It is in fact possible to archive many different timelines. While that might seem like a useless feature, it's often a lifesaver. Consider the situation where you aren't quite sure what point-in-time to recover to, and so have to do several point-in-time recoveries by trial and error until you find the best place to branch off from the old history. Without timelines this process would soon generate an unmanageable mess. With timelines, you can recover to *any* prior state, including states in timeline branches that you abandoned earlier.

Every time a new timeline is created, Postgres Pro creates a “timeline history” file that shows which timeline it branched off from and when. These history files are necessary to allow the system to pick the right WAL segment files when recovering from an archive that contains multiple timelines. Therefore, they are archived into the WAL archive area just like WAL segment files. The history files are just small text files, so it's cheap and appropriate to keep them around indefinitely (unlike the segment files which are large). You can, if you like, add comments to a history file to record your own notes about how and why this particular timeline was created. Such comments will be especially valuable when you have a thicket of different timelines as a result of experimentation.

The default behavior of recovery is to recover to the latest timeline found in the archive. If you wish to recover to the timeline that was current when the base backup was taken or into a specific child timeline (that is, you want to return to some state that was itself generated after a recovery attempt), you need to specify `current` or the target timeline ID in `recovery_target_timeline`. You cannot recover into timelines that branched off earlier than the base backup.

25.1.3. Tips and Examples

Some tips for configuring continuous archiving are given here.

25.1.3.1. Standalone Hot Backups

It is possible to use Postgres Pro's backup facilities to produce standalone hot backups. These are backups that cannot be used for point-in-time recovery, yet are typically much faster to backup and restore than `pg_dump` dumps. (They are also much larger than `pg_dump` dumps, so in some cases the speed advantage might be negated.)

25.1.3.2. Compressed Archive Logs

If archive storage size is a concern, you can use `gzip` to compress the archive files:

```
archive_command = 'gzip < %p > /mnt/server/archivedir/%f.gz'
```

You will then need to use `gunzip` during recovery:

```
restore_command = 'gunzip < /mnt/server/archivedir/%f.gz > %p'
```

25.1.3.3. `archive_command` Scripts

Many people choose to use scripts to define their `archive_command`, so that their `postgresql.conf` entry looks very simple:

```
archive_command = 'local_backup_script.sh "%p" "%f"'
```

Using a separate script file is advisable any time you want to use more than a single command in the archiving process. This allows all complexity to be managed within the script, which can be written in a popular scripting language such as `bash` or `perl`.

Examples of requirements that might be solved within a script include:

- Copying data to secure off-site data storage
- Batching WAL files so that they are transferred every three hours, rather than one at a time
- Interfacing with other backup and recovery software
- Interfacing with monitoring software to report errors

Tip

When using an `archive_command` script, it's desirable to enable [logging_collector](#). Any messages written to `stderr` from the script will then appear in the database server log, allowing complex configurations to be diagnosed easily if they fail.

25.1.4. Caveats

At this writing, there are several limitations of the continuous archiving technique. These will probably be fixed in future releases:

- If a `CREATE DATABASE` command is executed while a base backup is being taken, and then the template database that the `CREATE DATABASE` copied is modified while the base backup is still in progress, it is possible that recovery will cause those modifications to be propagated into the created database as well. This is of course undesirable. To avoid this risk, it is best not to modify any template databases while taking a base backup.
- `CREATE TABLESPACE` commands are WAL-logged with the literal absolute path, and will therefore be replayed as tablespace creations with the same absolute path. This might be undesirable if the WAL is being replayed on a different machine. It can be dangerous even if the WAL is being replayed on the same machine, but into a new data directory: the replay will still overwrite the contents of the original tablespace. To avoid potential gotchas of this sort, the best practice is to take a new base backup after creating or dropping tablespaces.

It should also be noted that the default WAL format is fairly bulky since it includes many disk page snapshots. These page snapshots are designed to support crash recovery, since we might need to fix partially-written disk pages. Depending on your system hardware and software, the risk of partial writes might be small enough to ignore, in which case you can significantly reduce the total volume of archived WAL files by turning off page snapshots using the [full_page_writes](#) parameter. (Read the notes and warnings in [Chapter 28](#) before you do so.) Turning off page snapshots does not prevent use of the WAL for PITR operations. An area for future development is to compress archived WAL data by removing unnecessary page copies even when `full_page_writes` is on. In the meantime, administrators might wish to reduce the number of page snapshots included in WAL by increasing the checkpoint interval parameters as much as feasible.

25.2. Backup and Recovery in Distributed System

This section describes basics of backup and recovery in Postgres Pro Shardman.

You can use the `probackup backup` command of the `shardmanctl` tool to perform a full binary consistent backup of a Postgres Pro Shardman cluster to the backup repository on the local host or S3-compatible object storage and the `probackup restore` command to perform a recovery from any backup from the repository. Full and partial (DELTA) backups are supported.

The Postgres Pro Shardman `pg_probackup` utility for creating consistent full and incremental backups was integrated into `shardman-utils`. `shardman-utils` uses the `pg_probackup` approach to store backups in a pre-created repository. In addition, the `pg_probackup` commands `archive-get` and `archive-push` are used to deliver WAL logs into the backup repository. Backup and restore modes use a passwordless SSH connection between the cluster nodes and the backup node.

Postgres Pro Shardman cluster configuration parameter `enable_csn_snapshot` must be set to `on`. This parameter is necessary for the cluster backup to be consistent. If this option is disabled, a consistent backup is not possible.

For consistent visibility of distributed transactions, the technique of global snapshots based on physical clocks is used. Similarly, it is possible to get a consistent snapshot for backups, only the time corresponding to the global snapshot must be mapped to the set of LSNs for each node. Such a set of consistent LSNs in a cluster is called a syncpoint. By getting the syncpoint and taking the LSN for each node in the cluster from it, we can make a backup of each node, which must necessarily contain that LSN. We can also recover to this LSN using the point in time recovery (PITR) mechanism.

The `probackup` command uses the `pg_probackup` utility and its options to create a cluster backup. In any case of using `probackup` commands for restoration, the node names, defined by the hostname or IP-address, must correspond to those that were in place at the time of the backup.

25.2.1. Requirements for Backup and Restore with `pg_probackup`

To backup and restore a Postgres Pro Shardman cluster via the `probackup` command, the following requirements must be met:

- On the backup host, Postgres Pro Shardman utilities must be installed into `/opt/pgpro/sdm-17/bin`.
- On the backup host and on each cluster node, `pg_probackup` must be installed into `/opt/pgpro/sdm-17/bin`.
- On the backup host, `postgres` Linux user and group must be created.
- Passwordless SSH connection between the backup host and each Postgres Pro Shardman cluster node for the `postgres` Linux user must be configured. To do this, on each node:
 - The `postgres` user must create the `.ssh` subdirectory in the `/var/lib/postgresql` directory and place there the keys required for the passwordless SSH connection.
 - To perform a backup/restore in a pretty large number of threads, such as 50 (`-j=50`, see [the section called “ backup ”](#) for details), `MaxSessions` and `MaxStartups` must be set to 100 for the backup host in the `/etc/ssh/sshd_config` file.

Note

Setting the number of threads (`-j` option) to a value greater than 10 for `shardmanctl probackup` may result in the actual number of SSH connections exceeding the maximum allowed number of simultaneous SSH connections on the backup host and consequently lead to an “ERROR: Agent error: kex_exchange_identification: Connection closed by remote host” error. To correct the error, either reduce the number of `probackup` threads or adjust the value of `MaxStartups` configuration parameter of the backup host. If SSH is set

up as a xinetd service on the backup host, adjust the value of the xinetd `per_source` configuration parameter rather than `MaxStartups`.

You can disable SSH for data copying by setting the `--storage-type` option to the `mount` or `S3` value (but SSH will be required to execute remote commands). Also this value will be automatically used in the restore process.

- A backup folder or bucket in the S3-compatible object storage must be created.
- Access for the `postgres` Linux user to the backup folder must be granted.
- `shardmanctl` utility must be run as `postgres` Linux user.
- `init` subcommand for the backup repository initialization must be successfully executed on the backup host.
- `archive-command add` subcommand for enabling `archive_command` for each replication group to stream WALs into the initialized repository must be successfully executed on the backup host.

25.2.2. pg_probackup Backup Process

`shardmanctl` conducts a backup task in several steps. The tool:

1. Takes necessary locks in `etcd` to prevent concurrent cluster-wide operations.
2. Connects to a random replication group and locks Postgres Pro Shardman metadata tables to prevent modification of foreign servers during the backup.
3. Dumps Postgres Pro Shardman metadata, stored in `etcd`, to a JSON file in the backup directory or bucket in the S3-compatible object storage.
4. To get backups from each replication group, concurrently runs `pg_probackup` using the configured `archive_command`.
5. Creates the syncpoint and gets LSNs for each replication group from the syncpoint data structure. Then uses the `pg_probackup archive-push` command to push WAL logs generated after finishing backup and the WAL file where syncpoint LSNs are present for each replication group.

For more information on the `backup` command, see [reference](#).

25.2.3. Cluster Restore from a Backup with pg_probackup

You can restore a backup on the same or compatible cluster. By *compatible* clusters, those that use the same Postgres Pro Shardman version and have the same number of replication groups are meant here.

`shardmanctl` can perform full restore, metadata-only restore, or schema-only restore. Metadata-only restore is useful if issues are encountered with the `etcd` instance, but DBMS data is not corrupted.

During metadata-only restore, `shardmanctl` restores `etcd` data from the dump created during the backup.

Important

Restoring metadata to an incompatible cluster can lead to catastrophic consequences, including data loss, since the metadata state can differ from the actual configuration layout. Do not perform metadata-only restore if there were cluster reconfigurations after the backup, such as addition or deletion of nodes, even if the same nodes were added back again.

Schema-only recovery restores only schema information without data. It can be useful if the scale of the data is large and the schema is needed for testing or checking.

During a full restore, `shardmanctl` checks whether the number of replication groups in the target cluster matches the number of replication groups in the backup. This means that you cannot restore on an empty

cluster, but need to add as many replication groups as necessary for the total number of them to match that of the cluster from which the backup was taken.

You can perform restoring only on the single shard using `--shard` parameter.

Also you can perform Point-in-Time Recovery using `--recovery-target-time` parameter. In this case, Postgres Pro Shardman finds the closest syncpoint to the specified timestamp and suggests restoring on that LSN. You can also specify a `--wal-limit` option to limit the number of WAL segments to be processed.

`shardmanctl` conducts full restore in several steps. The tool:

1. Takes the necessary locks in etcd to prevent concurrent cluster-wide operations and tries to assign replication groups in the backup to existing replication groups. If it cannot do this (for example, due to cluster incompatibility), the recovery fails.
2. Restores part of the etcd metadata: the cluster specification and parts of replication group definitions.
3. When the correct metadata is in place, runs `init` in PITR initialization mode with `RecoveryTargetName` set to the value of the syncpoint LSN from the backup info file. `DataRestoreCommand` and `RestoreCommand` are also taken from the backup info file. These commands are generated automatically during the backup phase, it is not recommended to make any corrections to the file containing the Postgres Pro Shardman cluster backup description. When restoring a cluster for each replication group, the WAL files containing the final LSN to restore will be requested automatically from the backup repository on the remote backup node via the `pg_probackup archive-get` command.
4. Waits for each replication group to recover.
5. Finally we need to enable `archive_command` back.

When performing a sequential restoration in Postgres Pro Shardman, be cautious of potential timeline conflicts within WAL (Write-Ahead Logging) segments. This issue commonly arises when restoring a database from a backup that was created at a certain point in time. If the database continues to operate and generate WAL segments after this backup, these new WAL segments are associated with a different timeline. During restoration, if the system tries to replay WAL segments from a different timeline - the one that diverged from the point of backup - it can lead to inconsistencies and conflicts. Additionally, after completing a restoration in Postgres Pro Shardman, it is strongly advised not to restore the database onto the same timeline or onto any timeline that precedes the one from which the backup was made.

For more information, see [reference](#).

25.2.4. Merging Backups with `pg_probackup`

The more incremental backups are created, the larger the total size of the backup catalog. To save the disk space, it is possible to merge incremental backups to their parent full backup by running the `merge` command with the backup ID of the most recent incremental backup specified:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://etcd2:2379,http://etcd3:2379
  probackup merge --backup-path backup_dir --backup-id backup_id
```

This command merges the backups that belong to a common incremental backup chain. If a full backup is specified, it is merged with its first incremental backup. If an incremental backup is specified, it is merged to its parent full backup, along with all the incremental backups between them. Once the merge is complete, the full backup covers all the merged data, and the incremental backups are removed as redundant. Thus, the merge operation is virtually the same as removing all the outdated backups from a full backup, but is much faster, especially for the large data volumes. It also saves I/O and network traffic when using `pg_probackup` in the remote mode.

Before merging, `pg_probackup` validates all the affected backups to ensure that they are valid. The current backup status can be seen by running the `show` command:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://etcd2:2379,http://etcd3:2379
  probackup show --backup-path backup_dir
```


For more information, see [reference](#).

25.2.5. Deleting Backups with `pg_probackup`

To delete a backup that is no longer needed, run the following command:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://etcd2:2379,http://etcd3:2379
  probackup delete --backup-path backup_dir --backup-id backup_id
```

This command deletes a backup with a specified `backup_id`, along with all the incremental backups that descend from this `backup_id`, if any. It allows you to delete some of the recent incremental backups without affecting the underlying full backup and other incremental backups that follow it.

To delete obsolete WAL files that are not needed for recovery, use the `--delete-wal` flag:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://etcd2:2379,http://etcd3:2379
  probackup delete --backup-path backup_dir --backup-id backup_id --delete-wal
```

For more information, see [reference](#).

Chapter 26. High Availability, Load Balancing, and Replication

Database servers can work together to allow a second server to take over quickly if the primary server fails (high availability), or to allow several computers to serve the same data (load balancing). Ideally, database servers could work together seamlessly. Web servers serving static web pages can be combined quite easily by merely load-balancing web requests to multiple machines. In fact, read-only database servers can be combined relatively easily too. Unfortunately, most database servers have a read/write mix of requests, and read/write servers are much harder to combine. This is because though read-only data needs to be placed on each server only once, a write to any server has to be propagated to all servers so that future read requests to those servers return consistent results.

This synchronization problem is the fundamental difficulty for servers working together. Because there is no single solution that eliminates the impact of the sync problem for all use cases, there are multiple solutions. Each solution addresses this problem in a different way, and minimizes its impact for a specific workload.

Some solutions deal with synchronization by allowing only one server to modify the data. Servers that can modify data are called read/write, *master* or *primary* servers. Servers that track changes in the primary are called *standby* or *secondary* servers. A standby server that cannot be connected to until it is promoted to a primary server is called a *warm standby* server, and one that can accept connections and serves read-only queries is called a *hot standby* server.

Some solutions are synchronous, meaning that a data-modifying transaction is not considered committed until all servers have committed the transaction. This guarantees that a failover will not lose any data and that all load-balanced servers will return consistent results no matter which server is queried. In contrast, asynchronous solutions allow some delay between the time of a commit and its propagation to the other servers, opening the possibility that some transactions might be lost in the switch to a backup server, and that load balanced servers might return slightly stale results. Asynchronous communication is used when synchronous would be too slow.

Solutions can also be categorized by their granularity. Some solutions can deal only with an entire database server, while others allow control at the per-table or per-database level.

Performance must be considered in any choice. There is usually a trade-off between functionality and performance. For example, a fully synchronous solution over a slow network might cut performance by more than half, while an asynchronous one might have a minimal performance impact.

The remainder of this section outlines various failover, replication, and load balancing solutions.

26.1. Comparison of Different Solutions

Shared Disk Failover

Shared disk failover avoids synchronization overhead by having only one copy of the database. It uses a single disk array that is shared by multiple servers. If the main database server fails, the standby server is able to mount and start the database as though it were recovering from a database crash. This allows rapid failover with no data loss.

Shared hardware functionality is common in network storage devices. Using a network file system is also possible, though care must be taken that the file system has full POSIX behavior (see [Section 18.2.2.1](#)). One significant limitation of this method is that if the shared disk array fails or becomes corrupt, the primary and standby servers are both nonfunctional. Another issue is that the standby server should never access the shared storage while the primary server is running.

File System (Block Device) Replication

A modified version of shared hardware functionality is file system replication, where all changes to a file system are mirrored to a file system residing on another computer. The only restriction is that

the mirroring must be done in a way that ensures the standby server has a consistent copy of the file system — specifically, writes to the standby must be done in the same order as those on the primary. DRBD is a popular file system replication solution for Linux.

Write-Ahead Log Shipping

Warm and hot standby servers can be kept current by reading a stream of write-ahead log (WAL) records. If the main server fails, the standby contains almost all of the data of the main server, and can be quickly made the new primary database server. This can be synchronous or asynchronous and can only be done for the entire database server.

A standby server can be implemented using file-based log shipping ([Section 26.2](#)) or streaming replication (see [Section 26.2.5](#)), or a combination of both. For information on hot standby, see [Section 26.4](#).

Logical Replication

Logical replication allows a database server to send a stream of data modifications to another server. Postgres Pro logical replication constructs a stream of logical data modifications from the WAL. Logical replication allows replication of data changes on a per-table basis. In addition, a server that is publishing its own changes can also subscribe to changes from another server, allowing data to flow in multiple directions. For more information on logical replication, see [Chapter 29](#). Through the logical decoding interface ([Chapter 50](#)), third-party extensions can also provide similar functionality.

Trigger-Based Primary-Standby Replication

A trigger-based replication setup typically funnels data modification queries to a designated primary server. Operating on a per-table basis, the primary server sends data changes (typically) asynchronously to the standby servers. Standby servers can answer queries while the primary is running, and may allow some local data changes or write activity. This form of replication is often used for offloading large analytical or data warehouse queries.

Slony-I is an example of this type of replication, with per-table granularity, and support for multiple standby servers. Because it updates the standby server asynchronously (in batches), there is possible data loss during fail over.

SQL-Based Replication Middleware

With SQL-based replication middleware, a program intercepts every SQL query and sends it to one or all servers. Each server operates independently. Read-write queries must be sent to all servers, so that every server receives any changes. But read-only queries can be sent to just one server, allowing the read workload to be distributed among them.

If queries are simply broadcast unmodified, functions like `random()`, `CURRENT_TIMESTAMP`, and sequences can have different values on different servers. This is because each server operates independently, and because SQL queries are broadcast rather than actual data changes. If this is unacceptable, either the middleware or the application must determine such values from a single source and then use those values in write queries. Care must also be taken that all transactions either commit or abort on all servers, perhaps using two-phase commit ([PREPARE TRANSACTION](#) and [COMMIT PREPARED](#)). Pgpool-II and Continuent Tungsten are examples of this type of replication.

Asynchronous Multimaster Replication

For servers that are not regularly connected or have slow communication links, like laptops or remote servers, keeping data consistent among servers is a challenge. Using asynchronous multimaster replication, each server works independently, and periodically communicates with the other servers to identify conflicting transactions. The conflicts can be resolved by users or conflict resolution rules. Bucardo is an example of this type of replication.

Synchronous Multimaster Replication

In synchronous multimaster replication, each server can accept write requests, and modified data is transmitted from the original server to every other server before each transaction commits. Heavy

write activity can cause excessive locking and commit delays, leading to poor performance. Read requests can be sent to any server. Some implementations use shared disk to reduce the communication overhead. Synchronous multimaster replication is best for mostly read workloads, though its big advantage is that any server can accept write requests — there is no need to partition workloads between primary and standby servers, and because the data changes are sent from one server to another, there is no problem with non-deterministic functions like `random()`.

Postgres Pro Enterprise provides `multimaster` extension that implements this type of replication. Using `multimaster`, you can configure a synchronous shared-nothing cluster that provides Online Transaction Processing (OLTP) scalability for read transactions, as well as ensures high availability with automatic disaster recovery.

Table 26.1 summarizes the capabilities of the various solutions listed above.

Table 26.1. High Availability, Load Balancing, and Replication Feature Matrix

Feature	Shared Disk	File System Repl.	Write-Ahead Log Shipping	Logical Repl.	Trigger-Based Repl.	SQL Repl. Middle-ware	Async. MM Repl.	Sync. MM Repl.	WAL Shipping and Control Channel
Popular examples	NAS	DRBD	built-in streaming repl.	built-in logical repl., pg-logical	Londiste, Slony	pgpool-II	Bucardo	multi-master extension of Postgres Pro Enterprise	biha extension of Postgres Pro Enterprise
Comm. method	shared disk	disk blocks	WAL	logical decoding	table rows	SQL	table rows	table rows and row locks	WAL and biha protocol
No special hardware required		•	•	•	•	•	•	•	•
Allows multiple primary servers				•		•	•	•	
No overhead on primary	•		•	•		•			•
No waiting for multiple servers	•		with sync off	with sync off	•		•		with sync off
Primary failure will never lose data	•	•	with sync on	with sync on		•		•	with sync on
Replicas accept read-only queries			with hot standby	•	•	•	•	•	•

Feature	Shared Disk	File System Repl.	Write-Ahead Log Shipping	Logical Repl.	Trigger-Based Repl.	SQL Repl. Middle-ware	Async. MM Repl.	Sync. MM Repl.	WAL Shipping and Control Channel
Per-table granularity				•	•		•	•	
No conflict resolution necessary	•	•	•		•	•		•	•
Automatic node failure detection, response, and cluster reconfiguration									•

There are a few solutions that do not fit into the above categories:

Data Partitioning

Data partitioning splits tables into data sets. Each set can be modified by only one server. For example, data can be partitioned by offices, e.g., London and Paris, with a server in each office. If queries combining London and Paris data are necessary, an application can query both servers, or primary/standby replication can be used to keep a read-only copy of the other office's data on each server.

Multiple-Server Parallel Query Execution

Many of the above solutions allow multiple servers to handle multiple queries, but none allow a single query to use multiple servers to complete faster. This solution allows multiple servers to work concurrently on a single query. It is usually accomplished by splitting the data among servers and having each server execute its part of the query and return results to a central server where they are combined and returned to the user. This can be implemented using the PL/Proxy tool set.

26.2. Log-Shipping Standby Servers

Continuous archiving can be used to create a *high availability* (HA) cluster configuration with one or more *standby servers* ready to take over operations if the primary server fails. This capability is widely referred to as *warm standby* or *log shipping*.

The primary and standby server work together to provide this capability, though the servers are only loosely coupled. The primary server operates in continuous archiving mode, while each standby server operates in continuous recovery mode, reading the WAL files from the primary. No changes to the database tables are required to enable this capability, so it offers low administration overhead compared to some other replication solutions. This configuration also has relatively low performance impact on the primary server.

Directly moving WAL records from one database server to another is typically described as log shipping. Postgres Pro implements file-based log shipping by transferring WAL records one file (WAL segment) at a time. WAL files (16MB) can be shipped easily and cheaply over any distance, whether it be to an adjacent system, another system at the same site, or another system on the far side of the globe. The bandwidth required for this technique varies according to the transaction rate of the primary server.

Record-based log shipping is more granular and streams WAL changes incrementally over a network connection (see [Section 26.2.5](#)).

It should be noted that log shipping is asynchronous, i.e., the WAL records are shipped after transaction commit. As a result, there is a window for data loss should the primary server suffer a catastrophic failure; transactions not yet shipped will be lost. The size of the data loss window in file-based log shipping can be limited by use of the `archive_timeout` parameter, which can be set as low as a few seconds. However such a low setting will substantially increase the bandwidth required for file shipping. Streaming replication (see [Section 26.2.5](#)) allows a much smaller window of data loss.

Recovery performance is sufficiently good that the standby will typically be only moments away from full availability once it has been activated. As a result, this is called a warm standby configuration which offers high availability. Restoring a server from an archived base backup and rollforward will take considerably longer, so that technique only offers a solution for disaster recovery, not high availability. A standby server can also be used for read-only queries, in which case it is called a *hot standby* server. See [Section 26.4](#) for more information.

26.2.1. Planning

It is usually wise to create the primary and standby servers so that they are as similar as possible, at least from the perspective of the database server. In particular, the path names associated with tablespaces will be passed across unmodified, so both primary and standby servers must have the same mount paths for tablespaces if that feature is used. Keep in mind that if `CREATE TABLESPACE` is executed on the primary, any new mount point needed for it must be created on the primary and all standby servers before the command is executed. Hardware need not be exactly the same, but experience shows that maintaining two identical systems is easier than maintaining two dissimilar ones over the lifetime of the application and system. In any case the hardware architecture must be the same — shipping from, say, a 32-bit to a 64-bit system will not work.

In general, log shipping between servers running different major Postgres Pro release levels is not possible. It is the policy of the Postgres Pro Global Development Group not to make changes to disk formats during minor release upgrades, so it is likely that running different minor release levels on primary and standby servers will work successfully. However, no formal support for that is offered and you are advised to keep primary and standby servers at the same release level as much as possible. When updating to a new minor release, the safest policy is to update the standby servers first — a new minor release is more likely to be able to read WAL files from a previous minor release than vice versa.

26.2.2. Standby Server Operation

A server enters standby mode if a `standby.signal` file exists in the data directory when the server is started.

In standby mode, the server continuously applies WAL received from the primary server. The standby server can read WAL from a WAL archive (see [restore_command](#)) or directly from the primary over a TCP connection (streaming replication). The standby server will also attempt to restore any WAL found in the standby cluster's `pg_wal` directory. That typically happens after a server restart, when the standby replays again WAL that was streamed from the primary before the restart, but you can also manually copy files to `pg_wal` at any time to have them replayed.

At startup, the standby begins by restoring all WAL available in the archive location, calling `restore_command`. Once it reaches the end of WAL available there and `restore_command` fails, it tries to restore any WAL available in the `pg_wal` directory. If that fails, and streaming replication has been configured, the standby tries to connect to the primary server and start streaming WAL from the last valid record found in archive or `pg_wal`. If that fails or streaming replication is not configured, or if the connection is later disconnected, the standby goes back to step 1 and tries to restore the file from the archive again. This loop of retries from the archive, `pg_wal`, and via streaming replication goes on until the server is stopped or is promoted.

Standby mode is exited and the server switches to normal operation when `pg_ctl promote` is run, or `pg_promote()` is called. Before failover, any WAL immediately available in the archive or in `pg_wal` will be restored, but no attempt is made to connect to the primary.

26.2.3. Preparing the Primary for Standby Servers

Set up continuous archiving on the primary to an archive directory accessible from the standby, as described in [Section 25.1](#). The archive location should be accessible from the standby even when the primary is down, i.e., it should reside on the standby server itself or another trusted server, not on the primary server.

If you want to use streaming replication, set up authentication on the primary server to allow replication connections from the standby server(s); that is, create a role and provide a suitable entry or entries in `pg_hba.conf` with the database field set to `replication`. Also ensure `max_wal_senders` is set to a sufficiently large value in the configuration file of the primary server. If replication slots will be used, ensure that `max_replication_slots` is set sufficiently high as well.

Take a base backup to bootstrap the standby server.

26.2.4. Setting Up a Standby Server

To set up the standby server, restore the base backup taken from primary server. Create a file `standby.signal` in the standby's cluster data directory. Set `restore_command` to a simple command to copy files from the WAL archive. If you plan to have multiple standby servers for high availability purposes, make sure that `recovery_target_timeline` is set to `latest` (the default), to make the standby server follow the timeline change that occurs at failover to another standby.

Note

`restore_command` should return immediately if the file does not exist; the server will retry the command again if necessary.

If you want to use streaming replication, fill in `primary_conninfo` with a libpq connection string, including the host name (or IP address) and any additional details needed to connect to the primary server. If the primary needs a password for authentication, the password needs to be specified in `primary_conninfo` as well.

If you're setting up the standby server for high availability purposes, set up WAL archiving, connections and authentication like the primary server, because the standby server will work as a primary server after failover.

If you're using a WAL archive, its size can be minimized using the `archive_cleanup_command` parameter to remove files that are no longer required by the standby server. The `pg_archivecleanup` utility is designed specifically to be used with `archive_cleanup_command` in typical single-standby configurations, see [pg_archivecleanup](#). Note however, that if you're using the archive for backup purposes, you need to retain files needed to recover from at least the latest base backup, even if they're no longer needed by the standby.

A simple example of configuration is:

```
primary_conninfo = 'host=192.168.1.50 port=5432 user=foo password=foopass options='-c wal_sender_timeout=5000''
restore_command = 'cp /path/to/archive/%f %p'
archive_cleanup_command = 'pg_archivecleanup /path/to/archive %r'
```

You can have any number of standby servers, but if you use streaming replication, make sure you set `max_wal_senders` high enough in the primary to allow them to be connected simultaneously.

26.2.5. Streaming Replication

Streaming replication allows a standby server to stay more up-to-date than is possible with file-based log shipping. The standby connects to the primary, which streams WAL records to the standby as they're generated, without waiting for the WAL file to be filled.

Streaming replication is asynchronous by default (see [Section 26.2.8](#)), in which case there is a small delay between committing a transaction in the primary and the changes becoming visible in the standby. This delay is however much smaller than with file-based log shipping, typically under one second assuming the standby is powerful enough to keep up with the load. With streaming replication, `archive_timeout` is not required to reduce the data loss window.

If you use streaming replication without file-based continuous archiving, the server might recycle old WAL segments before the standby has received them. If this occurs, the standby will need to be reinitialized from a new base backup. You can avoid this by setting `wal_keep_size` to a value large enough to ensure that WAL segments are not recycled too early, or by configuring a replication slot for the standby. If you set up a WAL archive that's accessible from the standby, these solutions are not required, since the standby can always use the archive to catch up provided it retains enough segments.

To use streaming replication, set up a file-based log-shipping standby server as described in [Section 26.2](#). The step that turns a file-based log-shipping standby into streaming replication standby is setting the `primary_conninfo` setting to point to the primary server. Set [listen_addresses](#) and authentication options (see `pg_hba.conf`) on the primary so that the standby server can connect to the replication pseudo-database on the primary server (see [Section 26.2.5.1](#)).

On systems that support the `keepalive` socket option, setting `tcp_keepalives_idle`, `tcp_keepalives_interval` and `tcp_keepalives_count` helps the primary promptly notice a broken connection.

Set the maximum number of concurrent connections from the standby servers (see [max_wal_senders](#) for details).

When the standby is started and `primary_conninfo` is set correctly, the standby will connect to the primary after replaying all WAL files available in the archive. If the connection is established successfully, you will see a `walreceiver` in the standby, and a corresponding `walsender` process in the primary.

26.2.5.1. Authentication

It is very important that the access privileges for replication be set up so that only trusted users can read the WAL stream, because it is easy to extract privileged information from it. Standby servers must authenticate to the primary as an account that has the `REPLICATION` privilege or a superuser. It is recommended to create a dedicated user account with `REPLICATION` and `LOGIN` privileges for replication. While `REPLICATION` privilege gives very high permissions, it does not allow the user to modify any data on the primary system, which the `SUPERUSER` privilege does.

Client authentication for replication is controlled by a `pg_hba.conf` record specifying replication in the `database` field. For example, if the standby is running on host IP `192.168.1.100` and the account name for replication is `foo`, the administrator can add the following line to the `pg_hba.conf` file on the primary:

```
# Allow the user "foo" from host 192.168.1.100 to connect to the primary
# as a replication standby if the user's password is correctly supplied.
#
# TYPE  DATABASE          USER              ADDRESS           METHOD
host    replication      foo               192.168.1.100/32 md5
```

The host name and port number of the primary, connection user name, and password are specified in the `primary_conninfo`. The password can also be set in the `~/.pgpass` file on the standby (specify replication in the `database` field). For example, if the primary is running on host IP `192.168.1.50`, port `5432`, the account name for replication is `foo`, and the password is `foopass`, the administrator can add the following line to the `postgresql.conf` file on the standby:

```
# The standby connects to the primary that is running on host 192.168.1.50
# and port 5432 as the user "foo" whose password is "foopass".
```



```
primary_conninfo = 'host=192.168.1.50 port=5432 user=foo password=foopass'
```

26.2.5.2. Monitoring

An important health indicator of streaming replication is the amount of WAL records generated in the primary, but not yet applied in the standby. You can calculate this lag by comparing the current WAL write location on the primary with the last WAL location received by the standby. These locations can be retrieved using `pg_current_wal_lsn` on the primary and `pg_last_wal_receive_lsn` on the standby, respectively (see [Table 9.95](#) and [Table 9.96](#) for details). The last WAL receive location in the standby is also displayed in the process status of the WAL receiver process, displayed using the `ps` command (see [Section 27.1](#) for details).

You can retrieve a list of WAL sender processes via the `pg_stat_replication` view. Large differences between `pg_current_wal_lsn` and the view's `sent_lsn` field might indicate that the primary server is under heavy load, while differences between `sent_lsn` and `pg_last_wal_receive_lsn` on the standby might indicate network delay, or that the standby is under heavy load.

On a hot standby, the status of the WAL receiver process can be retrieved via the `pg_stat_wal_receiver` view. A large difference between `pg_last_wal_replay_lsn` and the view's `flushed_lsn` indicates that WAL is being received faster than it can be replayed.

26.2.5.3. Repairing a Data Page from Standby

If you have [enabled checksums](#) when initializing your database, the I/O system can detect and report data corruption that may be caused by a system error or disk failure. On clusters with streaming replication configured between the primary and standby servers, a corrupted data page can be repaired automatically via streaming replication from standby.

To enable automatic page repair, set the `page_repair` variable to `true` on the primary server.

A typical page repair workflow is as follows:

1. Once the system detects a checksum mismatch when reading a data page, Postgres Pro Shardman blocks the backend and starts the repair process.
2. The WAL sender sends a request to the WAL receiver on the standby server to read the same data page.
3. If the standby server contains a valid page copy, the backend overwrites the corrupted page with this copy received from the WAL receiver and gets unblocked.

If the page on the standby is broken as well, an error occurs on standby, and the backend fails with the corresponding error message. In this case, you can try to restore the page from a backup.

Note

Page repair cannot be done in parallel. If another backend detects a different corrupted data page while page repair is in progress, it has to wait until the repair is complete.

26.2.6. Replication Slots

Replication slots provide an automated way to ensure that the primary server does not remove WAL segments until they have been received by all standbys, and that the primary does not remove rows which could cause a [recovery conflict](#) even when the standby is disconnected.

In lieu of using replication slots, it is possible to prevent the removal of old WAL segments using `wal_keep_size`, or by storing the segments in an archive using `archive_command` or `archive_library`. A disadvantage of these methods is that they often result in retaining more WAL segments than required, whereas replication slots retain only the number of segments known to be needed.

Similarly, [hot_standby_feedback](#) on its own, without also using a replication slot, provides protection against relevant rows being removed by vacuum, but provides no protection during any time period when the standby is not connected.

Caution

Beware that replication slots can cause the server to retain so many WAL segments that they fill up the space allocated for `pg_wal`. [max_slot_wal_keep_size](#) can be used to limit the size of WAL files retained by replication slots.

26.2.6.1. Querying and Manipulating Replication Slots

Each replication slot has a name, which can contain lower-case letters, numbers, and the underscore character.

Existing replication slots and their state can be seen in the [pg_replication_slots](#) view.

Slots can be created and dropped either via the streaming replication protocol (see [Section 56.4](#)) or via SQL functions (see [Section 9.28.6](#)).

26.2.6.2. Configuration Example

You can create a replication slot like this:

```
postgres=# SELECT * FROM pg_create_physical_replication_slot('node_a_slot');
 slot_name | lsn
-----+-----
 node_a_slot |
```

```
postgres=# SELECT slot_name, slot_type, active FROM pg_replication_slots;
 slot_name | slot_type | active
-----+-----+-----
 node_a_slot | physical | f
(1 row)
```

To configure the standby to use this slot, `primary_slot_name` should be configured on the standby. Here is a simple example:

```
primary_conninfo = 'host=192.168.1.50 port=5432 user=foo password=foopass'
primary_slot_name = 'node_a_slot'
```

26.2.7. Cascading Replication

The cascading replication feature allows a standby server to accept replication connections and stream WAL records to other standbys, acting as a relay. This can be used to reduce the number of direct connections to the primary and also to minimize inter-site bandwidth overheads.

A standby acting as both a receiver and a sender is known as a cascading standby. Standbys that are more directly connected to the primary are known as upstream servers, while those standby servers further away are downstream servers. Cascading replication does not place limits on the number or arrangement of downstream servers, though each standby connects to only one upstream server which eventually links to a single primary server.

A cascading standby sends not only WAL records received from the primary but also those restored from the archive. So even if the replication connection in some upstream connection is terminated, streaming replication continues downstream for as long as new WAL records are available.

Cascading replication is currently asynchronous. Synchronous replication (see [Section 26.2.8](#)) settings have no effect on cascading replication at present.

Hot standby feedback propagates upstream, whatever the cascaded arrangement.

If an upstream standby server is promoted to become the new primary, downstream servers will continue to stream from the new primary if `recovery_target_timeline` is set to 'latest' (the default).

To use cascading replication, set up the cascading standby so that it can accept replication connections (that is, set `max_wal_senders` and `hot_standby`, and configure `host-based authentication`). You will also need to set `primary_conninfo` in the downstream standby to point to the cascading standby.

26.2.8. Synchronous Replication

Postgres Pro streaming replication is asynchronous by default. If the primary server crashes then some transactions that were committed may not have been replicated to the standby server, causing data loss. The amount of data loss is proportional to the replication delay at the time of failover.

Synchronous replication offers the ability to confirm that all changes made by a transaction have been transferred to one or more synchronous standby servers. This extends that standard level of durability offered by a transaction commit. This level of protection is referred to as 2-safe replication in computer science theory, and group-1-safe (group-safe and 1-safe) when `synchronous_commit` is set to `remote_write`.

When requesting synchronous replication, each commit of a write transaction will wait until confirmation is received that the commit has been written to the write-ahead log on disk of both the primary and standby server. The only possibility that data can be lost is if both the primary and the standby suffer crashes at the same time. This can provide a much higher level of durability, though only if the sysadmin is cautious about the placement and management of the two servers. Waiting for confirmation increases the user's confidence that the changes will not be lost in the event of server crashes but it also necessarily increases the response time for the requesting transaction. The minimum wait time is the round-trip time between primary and standby.

Read-only transactions and transaction rollbacks need not wait for replies from standby servers. Sub-transaction commits do not wait for responses from standby servers, only top-level commits. Long running actions such as data loading or index building do not wait until the very final commit message. All two-phase commit actions require commit waits, including both prepare and commit.

A synchronous standby can be a physical replication standby or a logical replication subscriber. It can also be any other physical or logical WAL replication stream consumer that knows how to send the appropriate feedback messages. Besides the built-in physical and logical replication systems, this includes special programs such as `pg_receivewal` and `pg_recvlogical` as well as some third-party replication systems and custom programs. Check the respective documentation for details on synchronous replication support.

26.2.8.1. Basic Configuration

Once streaming replication has been configured, configuring synchronous replication requires only one additional configuration step: `synchronous_standby_names` must be set to a non-empty value. `synchronous_commit` must also be set to `on`, but since this is the default value, typically no change is required. (See [Section 19.5.1](#) and [Section 19.6.2](#).) This configuration will cause each commit to wait for confirmation that the standby has written the commit record to durable storage. `synchronous_commit` can be set by individual users, so it can be configured in the configuration file, for particular users or databases, or dynamically by applications, in order to control the durability guarantee on a per-transaction basis.

After a commit record has been written to disk on the primary, the WAL record is then sent to the standby. The standby sends reply messages each time a new batch of WAL data is written to disk, unless `wal_receiver_status_interval` is set to zero on the standby. In the case that `synchronous_commit` is set to `remote_apply`, the standby sends reply messages when the commit record is replayed, making the transaction visible. If the standby is chosen as a synchronous standby, according to the setting of `synchronous_standby_names` on the primary, the reply messages from that standby will be considered along with those from other synchronous standbys to decide when to release transactions waiting for confirmation that the commit record has been received. These parameters allow the administrator to

specify which standby servers should be synchronous standbys. Note that the configuration of synchronous replication is mainly on the primary. Named standbys must be directly connected to the primary; the primary knows nothing about downstream standby servers using cascaded replication.

Setting `synchronous_commit` to `remote_write` will cause each commit to wait for confirmation that the standby has received the commit record and written it out to its own operating system, but not for the data to be flushed to disk on the standby. This setting provides a weaker guarantee of durability than `on` does: the standby could lose the data in the event of an operating system crash, though not a PostgreSQL crash. However, it's a useful setting in practice because it can decrease the response time for the transaction. Data loss could only occur if both the primary and the standby crash and the database of the primary gets corrupted at the same time.

Setting `synchronous_commit` to `remote_apply` will cause each commit to wait until the current synchronous standbys report that they have replayed the transaction, making it visible to user queries. In simple cases, this allows for load balancing with causal consistency.

By default, if you add new synchronous standbys or reconnect a standby after a failure, the primary server is blocked until all standbys catch up with its current state. To minimize downtime, you can configure the `synchronous_standby_gap` setting on the primary. This option specifies the WAL data gap between the primary and its synchronous standbys, in kilobytes, which serves as a threshold for primary blocking. For the specified positive value to work, you should specify the `MIN` option along with other suitable options in the `synchronous_standby_names` parameter. For example:

```
synchronous_standby_names = 'FIRST 1 MIN 0 (s1) '
synchronous_standby_gap = 1
```

In this case the primary server will continue running while WAL data is being replayed on standby, until the specified gap is reached. Then the primary is blocked to complete the synchronization. The smaller the gap, the shorter the blocking time.

Users will stop waiting if a fast shutdown is requested. However, as when using asynchronous replication, the server will not fully shutdown until all outstanding WAL records are transferred to the currently connected standby servers.

26.2.8.2. Multiple Synchronous Standbys

Synchronous replication supports one or more synchronous standby servers; transactions will wait until all the standby servers which are considered as synchronous confirm receipt of their data. The number of synchronous standbys that transactions must wait for replies from is specified in `synchronous_standby_names`. This parameter also specifies a list of standby names and the method (`FIRST` and `ANY`) to choose synchronous standbys from the listed ones.

The method `FIRST` specifies a priority-based synchronous replication and makes transaction commits wait until their WAL records are replicated to the requested number of synchronous standbys chosen based on their priorities. The standbys whose names appear earlier in the list are given higher priority and will be considered as synchronous. Other standby servers appearing later in this list represent potential synchronous standbys. If any of the current synchronous standbys disconnects for whatever reason, it will be replaced immediately with the next-highest-priority standby.

An example of `synchronous_standby_names` for a priority-based multiple synchronous standbys is:

```
synchronous_standby_names = 'FIRST 2 (s1, s2, s3) '
```

In this example, if four standby servers `s1`, `s2`, `s3` and `s4` are running, the two standbys `s1` and `s2` will be chosen as synchronous standbys because their names appear early in the list of standby names. `s3` is a potential synchronous standby and will take over the role of synchronous standby when either of `s1` or `s2` fails. `s4` is an asynchronous standby since its name is not in the list.

The method `ANY` specifies a quorum-based synchronous replication and makes transaction commits wait until their WAL records are replicated to *at least* the requested number of synchronous standbys in the list.

An example of `synchronous_standby_names` for a quorum-based multiple synchronous standbys is:

```
synchronous_standby_names = 'ANY 2 (s1, s2, s3)'
```

In this example, if four standby servers `s1`, `s2`, `s3` and `s4` are running, transaction commits will wait for replies from at least any two standbys of `s1`, `s2` and `s3`. `s4` is an asynchronous standby since its name is not in the list.

The `FIRST` and `ANY` methods can be used together with the `MIN` option that specifies the minimum number of synchronous standbys that must be connected to the primary at all times. You can use this option if some of the synchronous standbys are unstable or have poor connection to the primary. It allows to prevent blocking of the primary server if one or more synchronous standbys temporarily go offline.

An example of using priority-based multiple synchronous standbys with the `MIN` option:

```
synchronous_standby_names = 'FIRST 3 MIN 2 (s1, s2, s3, s4)'
```

In this example, the primary will continue running even if only two standbys remain available.

For the quorum-based case, the syntax is quite similar:

```
synchronous_standby_names = 'ANY 3 MIN 2 (s1, s2, s3, s4)'
```

Once the connection is restored, the behavior of the primary server depends on the [synchronous_standby_gap](#) setting.

The synchronous states of standby servers can be viewed using the `pg_stat_replication` view.

26.2.8.3. Planning for Performance

Synchronous replication usually requires carefully planned and placed standby servers to ensure applications perform acceptably. Waiting doesn't utilize system resources, but transaction locks continue to be held until the transfer is confirmed. As a result, incautious use of synchronous replication will reduce performance for database applications because of increased response times and higher contention.

Postgres Pro allows the application developer to specify the durability level required via replication. This can be specified for the system overall, though it can also be specified for specific users or connections, or even individual transactions.

For example, an application workload might consist of: 10% of changes are important customer details, while 90% of changes are less important data that the business can more easily survive if it is lost, such as chat messages between users.

With synchronous replication options specified at the application level (on the primary) we can offer synchronous replication for the most important changes, without slowing down the bulk of the total workload. Application level options are an important and practical tool for allowing the benefits of synchronous replication for high performance applications.

You should consider that the network bandwidth must be higher than the rate of generation of WAL data.

26.2.8.4. Planning for High Availability

`synchronous_standby_names` specifies the number and names of synchronous standbys that transaction commits made when `synchronous_commit` is set to `on`, `remote_apply` or `remote_write` will wait for responses from. Such transaction commits may never be completed if any one of the synchronous standbys should crash.

The best solution for high availability is to ensure you keep as many synchronous standbys as requested. This can be achieved by naming multiple potential synchronous standbys using `synchronous_standby_names`.

In a priority-based synchronous replication, the standbys whose names appear earlier in the list will be used as synchronous standbys. Standbys listed after these will take over the role of synchronous standby if one of current ones should fail.

In a quorum-based synchronous replication, all the standbys appearing in the list will be used as candidates for synchronous standbys. Even if one of them should fail, the other standbys will keep performing the role of candidates of synchronous standby.

When a standby first attaches to the primary, it will not yet be properly synchronized. This is described as `catchup` mode. Once the lag between standby and primary reaches zero for the first time we move to real-time `streaming` state. The catch-up duration may be long immediately after the standby has been created. If the standby is shut down, then the catch-up period will increase according to the length of time the standby has been down. The standby is only able to become a synchronous standby once it has reached `streaming` state. This state can be viewed using the `pg_stat_replication` view.

If primary restarts while commits are waiting for acknowledgment, those waiting transactions will be marked fully committed once the primary database recovers. There is no way to be certain that all standbys have received all outstanding WAL data at time of the crash of the primary. Some transactions may not show as committed on the standby, even though they show as committed on the primary. The guarantee we offer is that the application will not receive explicit acknowledgment of the successful commit of a transaction until the WAL data is known to be safely received by all the synchronous standbys.

If you really cannot keep as many synchronous standbys as requested then you should decrease the number of synchronous standbys that transaction commits must wait for responses from in `synchronous_standby_names` (or disable it) and reload the configuration file on the primary server.

If the primary is isolated from remaining standby servers you should fail over to the best candidate of those other remaining standby servers.

If you need to re-create a standby server while transactions are waiting, make sure that the functions `pg_backup_start()` and `pg_backup_stop()` are run in a session with `synchronous_commit = off`, otherwise those requests will wait forever for the standby to appear.

26.2.9. Continuous Archiving in Standby

When continuous WAL archiving is used in a standby, there are two different scenarios: the WAL archive can be shared between the primary and the standby, or the standby can have its own WAL archive. When the standby has its own WAL archive, set `archive_mode` to `always`, and the standby will call the archive command for every WAL segment it receives, whether it's by restoring from the archive or by streaming replication. The shared archive can be handled similarly, but the `archive_command` or `archive_library` must test if the file being archived exists already, and if the existing file has identical contents. This requires more care in the `archive_command` or `archive_library`, as it must be careful to not overwrite an existing file with different contents, but return success if the exactly same file is archived twice. And all that must be done free of race conditions, if two servers attempt to archive the same file at the same time.

If `archive_mode` is set to `on`, the archiver is not enabled during recovery or standby mode. If the standby server is promoted, it will start archiving after the promotion, but will not archive any WAL or timeline history files that it did not generate itself. To get a complete series of WAL files in the archive, you must ensure that all WAL is archived, before it reaches the standby. This is inherently true with file-based log shipping, as the standby can only restore files that are found in the archive, but not if streaming replication is enabled. When a server is not in recovery mode, there is no difference between `on` and `always` modes.

26.3. Failover

If the primary server fails then the standby server should begin failover procedures.

If the standby server fails then no failover need take place. If the standby server can be restarted, even some time later, then the recovery process can also be restarted immediately, taking advantage of restartable recovery. If the standby server cannot be restarted, then a full new standby server instance should be created.

If the primary server fails and the standby server becomes the new primary, and then the old primary restarts, you must have a mechanism for informing the old primary that it is no longer the primary. This is sometimes known as STONITH (Shoot The Other Node In The Head), which is necessary to avoid situations where both systems think they are the primary, which will lead to confusion and ultimately data loss.

Many failover systems use just two systems, the primary and the standby, connected by some kind of heartbeat mechanism to continually verify the connectivity between the two and the viability of the primary. It is also possible to use a third system (called a witness server) to prevent some cases of inappropriate failover, but the additional complexity might not be worthwhile unless it is set up with sufficient care and rigorous testing.

Postgres Pro does not provide the system software required to identify a failure on the primary and notify the standby database server. Many such tools exist and are well integrated with the operating system facilities required for successful failover, such as IP address migration.

Once failover to the standby occurs, there is only a single server in operation. This is known as a degenerate state. The former standby is now the primary, but the former primary is down and might stay down. To return to normal operation, a standby server must be recreated, either on the former primary system when it comes up, or on a third, possibly new, system. The [pg_rewind](#) utility can be used to speed up this process on large clusters. Once complete, the primary and standby can be considered to have switched roles. Some people choose to use a third server to provide backup for the new primary until the new standby server is recreated, though clearly this complicates the system configuration and operational processes.

So, switching from primary to standby server can be fast but requires some time to re-prepare the failover cluster. Regular switching from primary to standby is useful, since it allows regular downtime on each system for maintenance. This also serves as a test of the failover mechanism to ensure that it will really work when you need it. Written administration procedures are advised.

If you have opted for logical replication slot synchronization (see [Section 50.2.3](#)), then before switching to the standby server, it is recommended to check if the logical slots synchronized on the standby server are ready for failover. This can be done by following the steps described in [Section 29.3](#).

To trigger failover of a log-shipping standby server, run `pg_ctl promote` or call `pg_promote()`. If you're setting up reporting servers that are only used to offload read-only queries from the primary, not for high availability purposes, you don't need to promote.

26.4. Hot Standby

Hot standby is the term used to describe the ability to connect to the server and run read-only queries while the server is in archive recovery or standby mode. This is useful both for replication purposes and for restoring a backup to a desired state with great precision. The term hot standby also refers to the ability of the server to move from recovery through to normal operation while users continue running queries and/or keep their connections open.

Running queries in hot standby mode is similar to normal query operation, though there are several usage and administrative differences explained below.

26.4.1. User's Overview

When the [hot standby](#) parameter is set to true on a standby server, it will begin accepting connections once the recovery has brought the system to a consistent state. All such connections are strictly read-only; not even temporary tables may be written.

The data on the standby takes some time to arrive from the primary server so there will be a measurable delay between primary and standby. Running the same query nearly simultaneously on both primary and standby might therefore return differing results. We say that data on the standby is *eventually consistent* with the primary. Once the commit record for a transaction is replayed on the standby, the changes

made by that transaction will be visible to any new snapshots taken on the standby. Snapshots may be taken at the start of each query or at the start of each transaction, depending on the current transaction isolation level. For more details, see [Section 13.2](#).

Transactions started during hot standby may issue the following commands:

- **Query access:** `SELECT`, `COPY TO`
- **Cursor commands:** `DECLARE`, `FETCH`, `CLOSE`
- **Settings:** `SHOW`, `SET`, `RESET`
- **Transaction management commands:**
 - `BEGIN`, `END`, `ABORT`, `START TRANSACTION`
 - `SAVEPOINT`, `RELEASE`, `ROLLBACK TO SAVEPOINT`
 - `EXCEPTION` blocks and other internal subtransactions
- `LOCK TABLE`, though only when explicitly in one of these modes: `ACCESS SHARE`, `ROW SHARE` or `ROW EXCLUSIVE`.
- **Plans and resources:** `PREPARE`, `EXECUTE`, `DEALLOCATE`, `DISCARD`
- **Plugins and extensions:** `LOAD`
- `UNLISTEN`

Transactions started during hot standby will never be assigned a transaction ID and cannot write to the system write-ahead log. Therefore, the following actions will produce error messages:

- **Data Manipulation Language (DML):** `INSERT`, `UPDATE`, `DELETE`, `MERGE`, `COPY FROM`, `TRUNCATE`. Note that there are no allowed actions that result in a trigger being executed during recovery. This restriction applies even to temporary tables, because table rows cannot be read or written without assigning a transaction ID, which is currently not possible in a hot standby environment.
- **Data Definition Language (DDL):** `CREATE`, `DROP`, `ALTER`, `COMMENT`. This restriction applies even to temporary tables, because carrying out these operations would require updating the system catalog tables.
- `SELECT ... FOR SHARE | UPDATE`, because row locks cannot be taken without updating the underlying data files.
- Rules on `SELECT` statements that generate DML commands.
- `LOCK` that explicitly requests a mode higher than `ROW EXCLUSIVE MODE`.
- `LOCK` in short default form, since it requests `ACCESS EXCLUSIVE MODE`.
- **Transaction management commands that explicitly set non-read-only state:**
 - `BEGIN READ WRITE`, `START TRANSACTION READ WRITE`
 - `SET TRANSACTION READ WRITE`, `SET SESSION CHARACTERISTICS AS TRANSACTION READ WRITE`
 - `SET transaction_read_only = off`
- **Two-phase commit commands:** `PREPARE TRANSACTION`, `COMMIT PREPARED`, `ROLLBACK PREPARED` because even read-only transactions need to write WAL in the prepare phase (the first phase of two phase commit).
- **Sequence updates:** `nextval()`, `setval()`
- `LISTEN`, `NOTIFY`

In normal operation, “read-only” transactions are allowed to use `LISTEN` and `NOTIFY`, so hot standby sessions operate under slightly tighter restrictions than ordinary read-only sessions. It is possible that some of these restrictions might be loosened in a future release.

During hot standby, the parameter `transaction_read_only` is always true and may not be changed. But as long as no attempt is made to modify the database, connections during hot standby will act much like any other database connection. If failover or switchover occurs, the database will switch to normal processing mode. Sessions will remain connected while the server changes mode. Once hot standby finishes, it will be possible to initiate read-write transactions (even from a session begun during hot standby).

Users can determine whether hot standby is currently active for their session by issuing `SHOW in_hot_standby`. (In server versions before 14, the `in_hot_standby` parameter did not exist; a workable substitute method for older servers is `SHOW transaction_read_only`.) In addition, a set of functions ([Table 9.96](#)) allow users to access information about the standby server. These allow you to write programs that are aware of the current state of the database. These can be used to monitor the progress of recovery, or to allow you to write complex programs that restore the database to particular states.

26.4.2. Handling Query Conflicts

The primary and standby servers are in many ways loosely connected. Actions on the primary will have an effect on the standby. As a result, there is potential for negative interactions or conflicts between them. The easiest conflict to understand is performance: if a huge data load is taking place on the primary then this will generate a similar stream of WAL records on the standby, so standby queries may contend for system resources, such as I/O.

There are also additional types of conflict that can occur with hot standby. These conflicts are *hard conflicts* in the sense that queries might need to be canceled and, in some cases, sessions disconnected to resolve them. The user is provided with several ways to handle these conflicts. Conflict cases include:

- Access Exclusive locks taken on the primary server, including both explicit `LOCK` commands and various DDL actions, conflict with table accesses in standby queries.
- Dropping a tablespace on the primary conflicts with standby queries using that tablespace for temporary work files.
- Dropping a database on the primary conflicts with sessions connected to that database on the standby.
- Application of a vacuum cleanup record from WAL conflicts with standby transactions whose snapshots can still “see” any of the rows to be removed.
- Application of a vacuum cleanup record from WAL conflicts with queries accessing the target page on the standby, whether or not the data to be removed is visible.

On the primary server, these cases simply result in waiting; and the user might choose to cancel either of the conflicting actions. However, on the standby there is no choice: the WAL-logged action already occurred on the primary so the standby must not fail to apply it. Furthermore, allowing WAL application to wait indefinitely may be very undesirable, because the standby's state will become increasingly far behind the primary's. Therefore, a mechanism is provided to forcibly cancel standby queries that conflict with to-be-applied WAL records.

An example of the problem situation is an administrator on the primary server running `DROP TABLE` on a table that is currently being queried on the standby server. Clearly the standby query cannot continue if the `DROP TABLE` is applied on the standby. If this situation occurred on the primary, the `DROP TABLE` would wait until the other query had finished. But when `DROP TABLE` is run on the primary, the primary doesn't have information about what queries are running on the standby, so it will not wait for any such standby queries. The WAL change records come through to the standby while the standby query is still running, causing a conflict. The standby server must either delay application of the WAL records (and everything after them, too) or else cancel the conflicting query so that the `DROP TABLE` can be applied.

When a conflicting query is short, it's typically desirable to allow it to complete by delaying WAL application for a little bit; but a long delay in WAL application is usually not desirable. So the cancel mechanism has parameters, `max_standby_archive_delay` and `max_standby_streaming_delay`, that define the maximum allowed delay in WAL application. Conflicting queries will be canceled once it has taken longer than the relevant delay setting to apply any newly-received WAL data. There are two parameters so that

different delay values can be specified for the case of reading WAL data from an archive (i.e., initial recovery from a base backup or “catching up” a standby server that has fallen far behind) versus reading WAL data via streaming replication.

In a standby server that exists primarily for high availability, it's best to set the delay parameters relatively short, so that the server cannot fall far behind the primary due to delays caused by standby queries. However, if the standby server is meant for executing long-running queries, then a high or even infinite delay value may be preferable. Keep in mind however that a long-running query could cause other sessions on the standby server to not see recent changes on the primary, if it delays application of WAL records.

Once the delay specified by `max_standby_archive_delay` or `max_standby_streaming_delay` has been exceeded, conflicting queries will be canceled. This usually results just in a cancellation error, although in the case of replaying a `DROP DATABASE` the entire conflicting session will be terminated. Also, if the conflict is over a lock held by an idle transaction, the conflicting session is terminated (this behavior might change in the future).

Canceled queries may be retried immediately (after beginning a new transaction, of course). Since query cancellation depends on the nature of the WAL records being replayed, a query that was canceled may well succeed if it is executed again.

Keep in mind that the delay parameters are compared to the elapsed time since the WAL data was received by the standby server. Thus, the grace period allowed to any one query on the standby is never more than the delay parameter, and could be considerably less if the standby has already fallen behind as a result of waiting for previous queries to complete, or as a result of being unable to keep up with a heavy update load.

The most common reason for conflict between standby queries and WAL replay is “early cleanup”. Normally, Postgres Pro allows cleanup of old row versions when there are no transactions that need to see them to ensure correct visibility of data according to MVCC rules. However, this rule can only be applied for transactions executing on the primary. So it is possible that cleanup on the primary will remove row versions that are still visible to a transaction on the standby.

Row version cleanup isn't the only potential cause of conflicts with standby queries. All index-only scans (including those that run on standbys) must use an MVCC snapshot that “agrees” with the visibility map. Conflicts are therefore required whenever `VACUUM` [sets a page as all-visible in the visibility map](#) containing one or more rows *not* visible to all standby queries. So even running `VACUUM` against a table with no updated or deleted rows requiring cleanup might lead to conflicts.

Users should be clear that tables that are regularly and heavily updated on the primary server will quickly cause cancellation of longer running queries on the standby. In such cases the setting of a finite value for `max_standby_archive_delay` or `max_standby_streaming_delay` can be considered similar to setting `statement_timeout`.

Remedial possibilities exist if the number of standby-query cancellations is found to be unacceptable. The first option is to set the parameter `hot_standby_feedback`, which prevents `VACUUM` from removing recently-dead rows and so cleanup conflicts do not occur. If you do this, you should note that this will delay cleanup of dead rows on the primary, which may result in undesirable table bloat. However, the cleanup situation will be no worse than if the standby queries were running directly on the primary server, and you are still getting the benefit of off-loading execution onto the standby. If standby servers connect and disconnect frequently, you might want to make adjustments to handle the period when `hot_standby_feedback` feedback is not being provided. For example, consider increasing `max_standby_archive_delay` so that queries are not rapidly canceled by conflicts in WAL archive files during disconnected periods. You should also consider increasing `max_standby_streaming_delay` to avoid rapid cancellations by newly-arrived streaming WAL entries after reconnection.

The number of query cancels and the reason for them can be viewed using the `pg_stat_database_conflicts` system view on the standby server. The `pg_stat_database` system view also contains summary information.

Users can control whether a log message is produced when WAL replay is waiting longer than `deadlock_timeout` for conflicts. This is controlled by the `log_recovery_conflict_waits` parameter.

26.4.3. Administrator's Overview

If `hot_standby` is on in `postgresql.conf` (the default value) and there is a `standby.signal` file present, the server will run in hot standby mode. However, it may take some time for hot standby connections to be allowed, because the server will not accept connections until it has completed sufficient recovery to provide a consistent state against which queries can run. During this period, clients that attempt to connect will be refused with an error message. To confirm the server has come up, either loop trying to connect from the application, or look for these messages in the server logs:

```
LOG:  entering standby mode
```

```
... then some time later ...
```

```
LOG:  consistent recovery state reached
```

```
LOG:  database system is ready to accept read-only connections
```

Consistency information is recorded once per checkpoint on the primary. It is not possible to enable hot standby when reading WAL written during a period when `wal_level` was not set to `replica` or `logical` on the primary. Reaching a consistent state can also be delayed in the presence of both of these conditions:

- A write transaction has more than 64 subtransactions
- Very long-lived write transactions

If you are running file-based log shipping ("warm standby"), you might need to wait until the next WAL file arrives, which could be as long as the `archive_timeout` setting on the primary.

The settings of some parameters determine the size of shared memory for tracking transaction IDs, locks, and prepared transactions. These shared memory structures must be no smaller on a standby than on the primary in order to ensure that the standby does not run out of shared memory during recovery. For example, if the primary had used a prepared transaction but the standby had not allocated any shared memory for tracking prepared transactions, then recovery could not continue until the standby's configuration is changed. The parameters affected are:

- `max_connections`
- `max_prepared_transactions`
- `max_locks_per_transaction`
- `max_wal_senders`
- `max_worker_processes`

The easiest way to ensure this does not become a problem is to have these parameters set on the standbys to values equal to or greater than on the primary. Therefore, if you want to increase these values, you should do so on all standby servers first, before applying the changes to the primary server. Conversely, if you want to decrease these values, you should do so on the primary server first, before applying the changes to all standby servers. Keep in mind that when a standby is promoted, it becomes the new reference for the required parameter settings for the standbys that follow it. Therefore, to avoid this becoming a problem during a switchover or failover, it is recommended to keep these settings the same on all standby servers.

The WAL tracks changes to these parameters on the primary. If a hot standby processes WAL that indicates that the current value on the primary is higher than its own value, it will log a warning and pause recovery, for example:

```
WARNING: hot standby is not possible because of insufficient parameter settings
```

```
DETAIL: max_connections = 80 is a lower setting than on the primary server, where its value was 100.
```

```
LOG: recovery has paused
DETAIL: If recovery is unpaused, the server will shut down.
HINT: You can then restart the server after making the necessary configuration changes.
```

At that point, the settings on the standby need to be updated and the instance restarted before recovery can continue. If the standby is not a hot standby, then when it encounters the incompatible parameter change, it will shut down immediately without pausing, since there is then no value in keeping it up.

It is important that the administrator select appropriate settings for `max_standby_archive_delay` and `max_standby_streaming_delay`. The best choices vary depending on business priorities. For example if the server is primarily tasked as a High Availability server, then you will want low delay settings, perhaps even zero, though that is a very aggressive setting. If the standby server is tasked as an additional server for decision support queries then it might be acceptable to set the maximum delay values to many hours, or even -1 which means wait forever for queries to complete.

Transaction status "hint bits" written on the primary are not WAL-logged, so data on the standby will likely re-write the hints again on the standby. Thus, the standby server will still perform disk writes even though all users are read-only; no changes occur to the data values themselves. Users will still write large sort temporary files and re-generate relcache info files, so no part of the database is truly read-only during hot standby mode. Note also that writes to remote databases using dblink module, and other operations outside the database using PL functions will still be possible, even though the transaction is read-only locally.

The following types of administration commands are not accepted during recovery mode:

- Data Definition Language (DDL): e.g., `CREATE INDEX`
- Privilege and Ownership: `GRANT`, `REVOKE`, `REASSIGN`
- Maintenance commands: `ANALYZE`, `VACUUM`, `CLUSTER`, `REINDEX`

Again, note that some of these commands are actually allowed during "read only" mode transactions on the primary.

As a result, you cannot create additional indexes that exist solely on the standby, nor statistics that exist solely on the standby. If these administration commands are needed, they should be executed on the primary, and eventually those changes will propagate to the standby.

`pg_cancel_backend()` and `pg_terminate_backend()` will work on user backends, but not the startup process, which performs recovery. `pg_stat_activity` does not show recovering transactions as active. As a result, `pg_prepared_xacts` is always empty during recovery. If you wish to resolve in-doubt prepared transactions, view `pg_prepared_xacts` on the primary and issue commands to resolve transactions there or resolve them after the end of recovery.

`pg_locks` will show locks held by backends, as normal. `pg_locks` also shows a virtual transaction managed by the startup process that owns all `AccessExclusiveLocks` held by transactions being replayed by recovery. Note that the startup process does not acquire locks to make database changes, and thus locks other than `AccessExclusiveLocks` do not show in `pg_locks` for the Startup process; they are just presumed to exist.

The Nagios plugin `check_pgsql` will work, because the simple information it checks for exists. The `check_postgres` monitoring script will also work, though some reported values could give different or confusing results. For example, last vacuum time will not be maintained, since no vacuum occurs on the standby. Vacuums running on the primary do still send their changes to the standby.

WAL file control commands will not work during recovery, e.g., `pg_backup_start`, `pg_switch_wal` etc.

Dynamically loadable modules work, including `pg_stat_statements`.

Advisory locks work normally in recovery, including deadlock detection. Note that advisory locks are never WAL logged, so it is impossible for an advisory lock on either the primary or the standby to conflict

with WAL replay. Nor is it possible to acquire an advisory lock on the primary and have it initiate a similar advisory lock on the standby. Advisory locks relate only to the server on which they are acquired.

Trigger-based replication systems such as Slony, Londiste and Bucardo won't run on the standby at all, though they will run happily on the primary server as long as the changes are not sent to standby servers to be applied. WAL replay is not trigger-based so you cannot relay from the standby to any system that requires additional database writes or relies on the use of triggers.

New OIDs cannot be assigned, though some UUID generators may still work as long as they do not rely on writing new status to the database.

Currently, temporary table creation is not allowed during read-only transactions, so in some cases existing scripts will not run correctly. This restriction might be relaxed in a later release. This is both an SQL standard compliance issue and a technical issue.

`DROP TABLESPACE` can only succeed if the tablespace is empty. Some standby users may be actively using the tablespace via their `temp_tablespaces` parameter. If there are temporary files in the tablespace, all active queries are canceled to ensure that temporary files are removed, so the tablespace can be removed and WAL replay can continue.

Running `DROP DATABASE` or `ALTER DATABASE ... SET TABLESPACE` on the primary will generate a WAL entry that will cause all users connected to that database on the standby to be forcibly disconnected. This action occurs immediately, whatever the setting of `max_standby_streaming_delay`. Note that `ALTER DATABASE ... RENAME` does not disconnect users, which in most cases will go unnoticed, though might in some cases cause a program confusion if it depends in some way upon database name.

In normal (non-recovery) mode, if you issue `DROP USER` or `DROP ROLE` for a role with login capability while that user is still connected then nothing happens to the connected user — they remain connected. The user cannot reconnect however. This behavior applies in recovery also, so a `DROP USER` on the primary does not disconnect that user on the standby.

The cumulative statistics system is active during recovery. All scans, reads, blocks, index usage, etc., will be recorded normally on the standby. However, WAL replay will not increment relation and database specific counters. I.e. replay will not increment `pg_stat_all_tables` columns (like `n_tup_ins`), nor will reads or writes performed by the startup process be tracked in the `pg_statio_` views, nor will associated `pg_stat_database` columns be incremented.

Autovacuum is not active during recovery. It will start normally at the end of recovery.

The checkpoint process and the background writer process are active during recovery. The checkpoint process will perform restartpoints (similar to checkpoints on the primary) and the background writer process will perform normal block cleaning activities. This can include updates of the hint bit information stored on the standby server. The `CHECKPOINT` command is accepted during recovery, though it performs a restartpoint rather than a new checkpoint.

26.4.4. Hot Standby Parameter Reference

Various parameters have been mentioned above in [Section 26.4.2](#) and [Section 26.4.3](#).

On the primary, the `wal_level` parameter can be used. `max_standby_archive_delay` and `max_standby_streaming_delay` have no effect if set on the primary.

On the standby, parameters `hot_standby`, `max_standby_archive_delay` and `max_standby_streaming_delay` can be used.

26.4.5. Caveats

There are several limitations of hot standby. These can and probably will be fixed in future releases:

- Full knowledge of running transactions is required before snapshots can be taken. Transactions that use large numbers of subtransactions (currently greater than 64) will delay the start of read-

only connections until the completion of the longest running write transaction. If this situation occurs, explanatory messages will be sent to the server log.

- Valid starting points for standby queries are generated at each checkpoint on the primary. If the standby is shut down while the primary is in a shutdown state, it might not be possible to re-enter hot standby until the primary is started up, so that it generates further starting points in the WAL logs. This situation isn't a problem in the most common situations where it might happen. Generally, if the primary is shut down and not available anymore, that's likely due to a serious failure that requires the standby being converted to operate as the new primary anyway. And in situations where the primary is being intentionally taken down, coordinating to make sure the standby becomes the new primary smoothly is also standard procedure.
- At the end of recovery, `AccessExclusiveLocks` held by prepared transactions will require twice the normal number of lock table entries. If you plan on running either a large number of concurrent prepared transactions that normally take `AccessExclusiveLocks`, or you plan on having one large transaction that takes many `AccessExclusiveLocks`, you are advised to select a larger value of `max_locks_per_transaction`, perhaps as much as twice the value of the parameter on the primary server. You need not consider this at all if your setting of `max_prepared_transactions` is 0.
- The Serializable transaction isolation level is not yet available in hot standby. (See [Section 13.2.3](#) and [Section 13.4.1](#) for details.) An attempt to set a transaction to the serializable isolation level in hot standby mode will generate an error.

Chapter 27. Monitoring Database Activity

A database administrator frequently wonders, “What is the system doing right now?” This chapter discusses how to find that out.

Several tools are available for monitoring database activity and analyzing performance. Most of this chapter is devoted to describing Postgres Pro Shardman's cumulative statistics system, but one should not neglect regular Unix monitoring programs such as `ps`, `top`, `iostat`, and `vmstat`. Also, once one has identified a poorly-performing query, further investigation might be needed using Postgres Pro Shardman's `EXPLAIN` command. [Section 14.1](#) discusses `EXPLAIN` and other methods for understanding the behavior of an individual query.

27.1. Standard Unix Tools

On most Unix platforms, Postgres Pro Shardman modifies its command title as reported by `ps`, so that individual server processes can readily be identified. A sample display is

```
$ ps auxww | grep ^postgres
postgres 15551 0.0 0.1 57536 7132 pts/0 S 18:02 0:00 postgres -i
postgres 15554 0.0 0.0 57536 1184 ? Ss 18:02 0:00 postgres: background
writer
postgres 15555 0.0 0.0 57536 916 ? Ss 18:02 0:00 postgres:
checkpointer
postgres 15556 0.0 0.0 57536 916 ? Ss 18:02 0:00 postgres: walwriter
postgres 15557 0.0 0.0 58504 2244 ? Ss 18:02 0:00 postgres: autovacuum
launcher
postgres 15582 0.0 0.0 58772 3080 ? Ss 18:04 0:00 postgres: joe runbug
127.0.0.1 idle
postgres 15606 0.0 0.0 58772 3052 ? Ss 18:07 0:00 postgres: tgl
regression [local] SELECT waiting
postgres 15610 0.0 0.0 58772 3056 ? Ss 18:07 0:00 postgres: tgl
regression [local] idle in transaction
```

(The appropriate invocation of `ps` varies across different platforms, as do the details of what is shown. This example is from a recent Linux system.) The first process listed here is the primary server process. The command arguments shown for it are the same ones used when it was launched. The next four processes are background worker processes automatically launched by the primary process. (The “autovacuum launcher” process will not be present if you have set the system not to run autovacuum.) Each of the remaining processes is a server process handling one client connection. Each such process sets its command line display in the form

```
postgres: user database host activity
```

The user, database, and (client) host items remain the same for the life of the client connection, but the activity indicator changes. The activity can be `idle` (i.e., waiting for a client command), `idle in transaction` (waiting for client inside a `BEGIN` block), or a command type name such as `SELECT`. Also, `waiting` is appended if the server process is presently waiting on a lock held by another session. In the above example we can infer that process 15606 is waiting for process 15610 to complete its transaction and thereby release some lock. (Process 15610 must be the blocker, because there is no other active session. In more complicated cases it would be necessary to look into the `pg_locks` system view to determine who is blocking whom.)

If `cluster_name` has been configured the cluster name will also be shown in `ps` output:

```
$ psql -c 'SHOW cluster_name'
cluster_name
-----
server1
(1 row)
```

```
$ ps aux|grep server1
postgres    27093  0.0  0.0  30096  2752 ?        Ss   11:34   0:00 postgres: server1:
background writer
...
```

If you have turned off [update_process_title](#) then the activity indicator is not updated; the process title is set only once when a new process is launched. On some platforms this saves a measurable amount of per-command overhead; on others it's insignificant.

Tip

Solaris requires special handling. You must use `/usr/ucb/ps`, rather than `/bin/ps`. You also must use two `w` flags, not just one. In addition, your original invocation of the `postgres` command must have a shorter `ps` status display than that provided by each server process. If you fail to do all three things, the `ps` output for each server process will be the original `postgres` command line.

27.2. The Cumulative Statistics System

Postgres Pro Shardman's *cumulative statistics system* supports collection and reporting of information about server activity. Presently, accesses to tables and indexes in both disk-block and individual-row terms are counted. The total number of rows in each table, and information about vacuum and analyze actions for each table are also counted. If enabled, calls to user-defined functions and the total time spent in each one are counted as well.

Postgres Pro Shardman also supports reporting dynamic information about exactly what is going on in the system right now, such as the exact command currently being executed by other server processes, and which other connections exist in the system. This facility is independent of the cumulative statistics system.

27.2.1. Statistics Collection Configuration

Since collection of statistics adds some overhead to query execution, the system can be configured to collect or not collect information. This is controlled by configuration parameters that are normally set in `postgresql.conf`. (See [Chapter 19](#) for details about setting configuration parameters.)

The parameter [track_activities](#) enables monitoring of the current command being executed by any server process.

The parameter [track_counts](#) controls whether cumulative statistics are collected about table and index accesses.

The parameter [track_functions](#) enables tracking of usage of user-defined functions.

The parameter [track_io_timing](#) enables monitoring of block read, write, extend, and fsync times.

The parameter [track_wal_io_timing](#) enables monitoring of WAL write and fsync times.

Normally these parameters are set in `postgresql.conf` so that they apply to all server processes, but it is possible to turn them on or off in individual sessions using the [SET](#) command. (To prevent ordinary users from hiding their activity from the administrator, only superusers are allowed to change these parameters with `SET`.)

Cumulative statistics are collected in shared memory. Every Postgres Pro Shardman process collects statistics locally, then updates the shared data at appropriate intervals. When a server, including a physical replica, shuts down cleanly, a permanent copy of the statistics data is stored in the `pg_stat` subdirectory, so that statistics can be retained across server restarts. In contrast, when starting from an unclean shutdown (e.g., after an immediate shutdown, a server crash, starting from a base backup, and point-in-time recovery), all statistics counters are reset.

27.2.2. Viewing Statistics

Several predefined views, listed in [Table 27.1](#), are available to show the current state of the system. There are also several other views, listed in [Table 27.2](#), available to show the accumulated statistics. Alternatively, one can build custom views using the underlying cumulative statistics functions, as discussed in [Section 27.2.28](#).

When using the cumulative statistics views and functions to monitor collected data, it is important to realize that the information does not update instantaneously. Each individual server process flushes out accumulated statistics to shared memory just before going idle, but not more frequently than once per `PGSTAT_MIN_INTERVAL` milliseconds (1 second unless altered while building the server); so a query or transaction still in progress does not affect the displayed totals and the displayed information lags behind actual activity. However, current-query information collected by `track_activities` is always up-to-date.

Another important point is that when a server process is asked to display any of the accumulated statistics, accessed values are cached until the end of its current transaction in the default configuration. So the statistics will show static information as long as you continue the current transaction. Similarly, information about the current queries of all sessions is collected when any such information is first requested within a transaction, and the same information will be displayed throughout the transaction. This is a feature, not a bug, because it allows you to perform several queries on the statistics and correlate the results without worrying that the numbers are changing underneath you. When analyzing statistics interactively, or with expensive queries, the time delta between accesses to individual statistics can lead to significant skew in the cached statistics. To minimize skew, `stats_fetch_consistency` can be set to `snapshot`, at the price of increased memory usage for caching not-needed statistics data. Conversely, if it's known that statistics are only accessed once, caching accessed statistics is unnecessary and can be avoided by setting `stats_fetch_consistency` to `none`. You can invoke `pg_stat_clear_snapshot()` to discard the current transaction's statistics snapshot or cached values (if any). The next use of statistical information will (when in snapshot mode) cause a new snapshot to be built or (when in cache mode) accessed statistics to be cached.

A transaction can also see its own statistics (not yet flushed out to the shared memory statistics) in the views `pg_stat_xact_all_tables`, `pg_stat_xact_sys_tables`, `pg_stat_xact_user_tables`, and `pg_stat_xact_user_functions`. These numbers do not act as stated above; instead they update continuously throughout the transaction.

Some of the information in the dynamic statistics views shown in [Table 27.1](#) is security restricted. Ordinary users can only see all the information about their own sessions (sessions belonging to a role that they are a member of). In rows about other sessions, many columns will be null. Note, however, that the existence of a session and its general properties such as its sessions user and database are visible to all users. Superusers and roles with privileges of built-in role `pg_read_all_stats` (see also [Section 21.5](#)) can see all the information about all sessions.

Table 27.1. Dynamic Statistics Views

View Name	Description
<code>pg_stat_activity</code>	One row per server process, showing information related to the current activity of that process, such as state and current query. See pg_stat_activity for details.
<code>pgpro_stat_wal_activity</code>	One row per server process, showing the size of WAL files generated by this process, as well as all the data provided in the <code>pg_stat_activity</code> view. See pgpro_stat_wal_activity for details.
<code>pg_stat_replication</code>	One row per WAL sender process, showing statistics about replication to that sender's connected

View Name	Description
	standby server. See pg_stat_replication for details.
<code>pg_stat_wal_receiver</code>	Only one row, showing statistics about the WAL receiver from that receiver's connected server. See pg_stat_wal_receiver for details.
<code>pg_stat_recovery_prefetch</code>	Only one row, showing statistics about blocks prefetched during recovery. See pg_stat_recovery_prefetch for details.
<code>pg_stat_subscription</code>	At least one row per subscription, showing information about the subscription workers. See pg_stat_subscription for details.
<code>pg_stat_ssl</code>	One row per connection (regular and replication), showing information about SSL used on this connection. See pg_stat_ssl for details.
<code>pg_stat_gssapi</code>	One row per connection (regular and replication), showing information about GSSAPI authentication and encryption used on this connection. See pg_stat_gssapi for details.
<code>pg_stat_progress_analyze</code>	One row for each backend (including autovacuum worker processes) running <code>ANALYZE</code> , showing current progress. See Section 27.4.1 .
<code>pg_stat_progress_create_index</code>	One row for each backend running <code>CREATE INDEX</code> or <code>REINDEX</code> , showing current progress. See Section 27.4.4 .
<code>pg_stat_progress_vacuum</code>	One row for each backend (including autovacuum worker processes) running <code>VACUUM</code> , showing current progress. See Section 27.4.5 .
<code>pg_stat_progress_cluster</code>	One row for each backend running <code>CLUSTER</code> or <code>VACUUM FULL</code> , showing current progress. See Section 27.4.2 .
<code>pg_stat_progress_basebackup</code>	One row for each WAL sender process streaming a base backup, showing current progress. See Section 27.4.6 .
<code>pg_stat_progress_copy</code>	One row for each backend running <code>COPY</code> , showing current progress. See Section 27.4.3 .

Table 27.2. Collected Statistics Views

View Name	Description
<code>pg_stat_archiver</code>	One row only, showing statistics about the WAL archiver process's activity. See pg_stat_archiver for details.
<code>pg_stat_bgwriter</code>	One row only, showing statistics about the background writer process's activity. See pg_stat_bgwriter for details.
<code>pg_stat_checkpoint</code>	One row only, showing statistics about the checkpoint process's activity. See pg_stat_checkpoint for details.
<code>pg_stat_database</code>	One row per database, showing database-wide statistics. See pg_stat_database for details.

View Name	Description
<code>pg_stat_database_conflicts</code>	One row per database, showing database-wide statistics about query cancels due to conflict with recovery on standby servers. See pg_stat_database_conflicts for details.
<code>pg_stat_io</code>	One row for each combination of backend type, context, and target object containing cluster-wide I/O statistics. See pg_stat_io for details.
<code>pg_stat_replication_slots</code>	One row per replication slot, showing statistics about the replication slot's usage. See pg_stat_replication_slots for details.
<code>pg_stat_slru</code>	One row per SLRU, showing statistics of operations. See pg_stat_slru for details.
<code>pg_stat_subscription_stats</code>	One row per subscription, showing statistics about errors. See pg_stat_subscription_stats for details.
<code>pg_stat_wal</code>	One row only, showing statistics about WAL activity. See pg_stat_wal for details.
<code>pg_stat_all_tables</code>	One row for each table in the current database, showing statistics about accesses to that specific table. See pg_stat_all_tables for details.
<code>pg_stat_sys_tables</code>	Same as <code>pg_stat_all_tables</code> , except that only system tables are shown.
<code>pg_stat_user_tables</code>	Same as <code>pg_stat_all_tables</code> , except that only user tables are shown.
<code>pg_stat_xact_all_tables</code>	Similar to <code>pg_stat_all_tables</code> , but counts actions taken so far within the current transaction (which are <i>not</i> yet included in <code>pg_stat_all_tables</code> and related views). The columns for numbers of live and dead rows and vacuum and analyze actions are not present in this view.
<code>pg_stat_xact_sys_tables</code>	Same as <code>pg_stat_xact_all_tables</code> , except that only system tables are shown.
<code>pg_stat_xact_user_tables</code>	Same as <code>pg_stat_xact_all_tables</code> , except that only user tables are shown.
<code>pg_stat_all_indexes</code>	One row for each index in the current database, showing statistics about accesses to that specific index. See pg_stat_all_indexes for details.
<code>pg_stat_sys_indexes</code>	Same as <code>pg_stat_all_indexes</code> , except that only indexes on system tables are shown.
<code>pg_stat_user_indexes</code>	Same as <code>pg_stat_all_indexes</code> , except that only indexes on user tables are shown.
<code>pg_stat_user_functions</code>	One row for each tracked function, showing statistics about executions of that function. See pg_stat_user_functions for details.
<code>pg_stat_xact_user_functions</code>	Similar to <code>pg_stat_user_functions</code> , but counts only calls during the current transaction (which are <i>not</i> yet included in <code>pg_stat_user_functions</code>).

View Name	Description
<code>pg_statio_all_tables</code>	One row for each table in the current database, showing statistics about I/O on that specific table. See pg_statio_all_tables for details.
<code>pg_statio_sys_tables</code>	Same as <code>pg_statio_all_tables</code> , except that only system tables are shown.
<code>pg_statio_user_tables</code>	Same as <code>pg_statio_all_tables</code> , except that only user tables are shown.
<code>pg_statio_all_indexes</code>	One row for each index in the current database, showing statistics about I/O on that specific index. See pg_statio_all_indexes for details.
<code>pg_statio_sys_indexes</code>	Same as <code>pg_statio_all_indexes</code> , except that only indexes on system tables are shown.
<code>pg_statio_user_indexes</code>	Same as <code>pg_statio_all_indexes</code> , except that only indexes on user tables are shown.
<code>pg_statio_all_sequences</code>	One row for each sequence in the current database, showing statistics about I/O on that specific sequence. See pg_statio_all_sequences for details.
<code>pg_statio_sys_sequences</code>	Same as <code>pg_statio_all_sequences</code> , except that only system sequences are shown. (Presently, no system sequences are defined, so this view is always empty.)
<code>pg_statio_user_sequences</code>	Same as <code>pg_statio_all_sequences</code> , except that only user sequences are shown.

The per-index statistics are particularly useful to determine which indexes are being used and how effective they are.

The `pg_stat_io` and `pg_statio_set` of views are useful for determining the effectiveness of the buffer cache. They can be used to calculate a cache hit ratio. Note that while Postgres Pro Shardman's I/O statistics capture most instances in which the kernel was invoked in order to perform I/O, they do not differentiate between data which had to be fetched from disk and that which already resided in the kernel page cache. Users are advised to use the Postgres Pro Shardman statistics views in combination with operating system utilities for a more complete picture of their database's I/O performance.

27.2.3. `pg_stat_activity`

The `pg_stat_activity` view will have one row per server process, showing information related to the current activity of that process.

Table 27.3. `pg_stat_activity` View

Column	Type	Description
<code>datid</code>	<code>oid</code>	OID of the database this backend is connected to
<code>datname</code>	<code>name</code>	Name of the database this backend is connected to
<code>pid</code>	<code>integer</code>	Process ID of this backend
<code>leader_pid</code>	<code>integer</code>	Process ID of the parallel group leader if this process is a parallel query worker, or process ID of the leader apply worker if this process is a parallel apply worker. <code>NULL</code> indicates that

Column	Type	Description
		this process is a parallel group leader or leader apply worker, or does not participate in any parallel operation.
usesysid	oid	OID of the user logged into this backend
username	name	Name of the user logged into this backend
application_name	text	Name of the application that is connected to this backend
client_addr	inet	IP address of the client connected to this backend. If this field is null, it indicates either that the client is connected via a Unix socket on the server machine or that this is an internal process such as autovacuum.
client_hostname	text	Host name of the connected client, as reported by a reverse DNS lookup of <code>client_addr</code> . This field will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	integer	TCP port number that the client is using for communication with this backend, or -1 if a Unix socket is used. If this field is null, it indicates that this is an internal server process.
backend_start	timestamp with time zone	Time when this process was started. For client backends, this is the time the client connected to the server.
xact_start	timestamp with time zone	Time when this process' current transaction was started, or null if no transaction is active. If the current query is the first of its transaction, this column is equal to the <code>query_start</code> column.
query_start	timestamp with time zone	Time when the currently active query was started, or if <code>state</code> is not active, when the last query was started
state_change	timestamp with time zone	Time when the <code>state</code> was last changed
wait_event_type	text	The type of event for which the backend is waiting, if any; otherwise NULL. See Table 27.4 .
wait_event	text	Wait event name if backend is currently waiting, otherwise NULL. See Table 27.5 through Table 27.13 .
state	text	Current overall state of this backend. Possible values are: <ul style="list-style-type: none"> <code>active</code>: The backend is executing a query. <code>idle</code>: The backend is waiting for a new client command. <code>idle in transaction</code>: The backend is in a transaction, but is not currently executing a query. <code>idle in transaction (aborted)</code>: This state is similar to <code>idle in transaction</code>, except one of the statements in the transaction caused an error. <code>fastpath function call</code>: The backend is executing a fast-path function. <code>disabled</code>: This state is reported if track_activities is disabled in this backend.
backend_xid	xid	

Column Type	Description
	Top-level transaction identifier of this backend, if any; see Section 67.1 .
<code>backend_xmin</code> <code>xid</code>	The current backend's <code>xmin</code> horizon.
<code>query_id</code> <code>bigint</code>	Identifier of this backend's most recent query. If <code>state</code> is <code>active</code> this field shows the identifier of the currently executing query. In all other states, it shows the identifier of last query that was executed. Query identifiers are not computed by default so this field will be null unless compute_query_id parameter is enabled or a third-party module that computes query identifiers is configured.
<code>query</code> <code>text</code>	Text of this backend's most recent query. If <code>state</code> is <code>active</code> this field shows the currently executing query. In all other states, it shows the last query that was executed. By default the query text is truncated at 1024 bytes; this value can be changed via the parameter track_activity_query_size .
<code>backend_type</code> <code>text</code>	Type of current backend. Possible types are <code>autovacuum launcher</code> , <code>autovacuum worker</code> , <code>logical replication launcher</code> , <code>logical replication worker</code> , <code>parallel worker</code> , <code>background writer</code> , <code>client backend</code> , <code>checkpointer</code> , <code>archiver</code> , <code>standalone backend</code> , <code>startup</code> , <code>walreceiver</code> , <code>walsender</code> , <code>walwriter</code> and <code>walsummarizer</code> . In addition, background workers registered by extensions may have additional types.

Note

The `wait_event` and `state` columns are independent. If a backend is in the `active` state, it may or may not be waiting on some event. If the state is `active` and `wait_event` is non-null, it means that a query is being executed, but is being blocked somewhere in the system.

Table 27.4. Wait Event Types

Wait Event Type	Description
Activity	The server process is idle. This event type indicates a process waiting for activity in its main processing loop. <code>wait_event</code> will identify the specific wait point; see Table 27.5 .
BufferPin	The server process is waiting for exclusive access to a data buffer. Buffer pin waits can be protracted if another process holds an open cursor that last read data from the buffer in question. See Table 27.6 .
Client	The server process is waiting for activity on a socket connected to a user application. Thus, the server expects something to happen that is independent of its internal processes. <code>wait_event</code> will identify the specific wait point; see Table 27.7 .
Extension	The server process is waiting for some condition defined by an extension module. See Table 27.8 .
InjectionPoint	The server process is waiting for an injection point to reach an outcome defined in a test. See Section 39.10.13 for more details. This type has no predefined wait points.

Wait Event Type	Description
IO	The server process is waiting for an I/O operation to complete. <code>wait_event</code> will identify the specific wait point; see Table 27.9 .
IPC	The server process is waiting for some interaction with another server process. <code>wait_event</code> will identify the specific wait point; see Table 27.10 .
Lock	The server process is waiting for a heavyweight lock. Heavyweight locks, also known as lock manager locks or simply locks, primarily protect SQL-visible objects such as tables. However, they are also used to ensure mutual exclusion for certain internal operations such as relation extension. <code>wait_event</code> will identify the type of lock awaited; see Table 27.11 .
LWLock	The server process is waiting for a lightweight lock. Most such locks protect a particular data structure in shared memory. <code>wait_event</code> will contain a name identifying the purpose of the lightweight lock. (Some locks have specific names; others are part of a group of locks each with a similar purpose.) See Table 27.12 .
Timeout	The server process is waiting for a timeout to expire. <code>wait_event</code> will identify the specific wait point; see Table 27.13 .

Table 27.5. Wait Events of Type Activity

Activity Wait Event	Description
ArchiverMain	Waiting in main loop of archiver process.
AutovacuumMain	Waiting in main loop of autovacuum launcher process.
BgfreezerMain	Waiting for background freezer to be enabled.
BgwriterHibernate	Waiting in background writer process, hibernating.
BgwriterMain	Waiting in main loop of background writer process.
CfsGcEnable	Waiting for CFS garbage collection to be enabled.
CheckpointMain	Waiting in main loop of checkpoint process.
LogicalApplyMain	Waiting in main loop of logical replication apply process.
LogicalLauncherMain	Waiting in main loop of logical replication launcher process.
LogicalParallelApplyMain	Waiting in main loop of logical replication parallel apply process.
PriorityMain	Waiting in main loop of priority worker.
RecoveryWalStream	Waiting in main loop of startup process for WAL to arrive, during streaming recovery.
ReplicationRepairMain	Waiting in main loop of replication repair worker.
ReplicationSlotsyncMain	Waiting in main loop of slot sync worker.

Activity Wait Event	Description
ReplicationSlotsyncShutdown	Waiting for slot sync worker to shut down.
SysloggerMain	Waiting in main loop of syslogger process.
WalReceiverMain	Waiting in main loop of WAL receiver process.
WalSenderMain	Waiting in main loop of WAL sender process.
WalSummarizerWal	Waiting in WAL summarizer for more WAL to be generated.
WalWriterMain	Waiting in main loop of WAL writer process.

Table 27.6. Wait Events of Type Bufferpin

BufferPin Wait Event	Description
BufferPin	Waiting to acquire an exclusive pin on a buffer.

Table 27.7. Wait Events of Type Client

Client Wait Event	Description
ClientRead	Waiting to read data from the client.
ClientWrite	Waiting to write data to the client.
GssOpenServer	Waiting to read data from the client while establishing a GSSAPI session.
LibpqwalreceiverConnect	Waiting in WAL receiver to establish connection to remote server.
LibpqwalreceiverReceive	Waiting in WAL receiver to receive data from remote server.
SslOpenServer	Waiting for SSL while attempting connection.
WaitForStandbyConfirmation	Waiting for WAL to be received and flushed by the physical standby.
WalSenderWaitForWal	Waiting for WAL to be flushed in WAL sender process.
WalSenderWriteData	Waiting for any activity when processing replies from WAL receiver in WAL sender process.

Table 27.8. Wait Events of Type Extension

Extension Wait Event	Description
Extension	Waiting in an extension.

Table 27.9. Wait Events of Type Io

IO Wait Event	Description
BasebackupRead	Waiting for base backup to read from a file.
BasebackupSync	Waiting for data written by a base backup to reach durable storage.
BasebackupWrite	Waiting for base backup to write to a file.
BuffblockRepair	Waiting for repair of the corrupted block using replication protocol.
BufffileRead	Waiting for a read from a buffered file.
BufffileTruncate	Waiting for a buffered file to be truncated.
BufffileWrite	Waiting for a write to a buffered file.
ControlFileRead	Waiting for a read from the pg_control file.

IO Wait Event	Description
ControlFileSync	Waiting for the <code>pg_control</code> file to reach durable storage.
ControlFileSyncUpdate	Waiting for an update to the <code>pg_control</code> file to reach durable storage.
ControlFileWrite	Waiting for a write to the <code>pg_control</code> file.
ControlFileWriteUpdate	Waiting for a write to update the <code>pg_control</code> file.
CopyFileRead	Waiting for a read during a file copy operation.
CopyFileWrite	Waiting for a write during a file copy operation.
DataFileExtend	Waiting for a relation data file to be extended.
DataFileFlush	Waiting for a relation data file to reach durable storage.
DataFileImmediateSync	Waiting for an immediate synchronization of a relation data file to durable storage.
DataFilePrefetch	Waiting for an asynchronous prefetch from a relation data file.
DataFileRead	Waiting for a read from a relation data file.
DataFileSync	Waiting for changes to a relation data file to reach durable storage.
DataFileTruncate	Waiting for a relation data file to be truncated.
DataFileWrite	Waiting for a write to a relation data file.
DsmAllocate	Waiting for a dynamic shared memory segment to be allocated.
DsmFillZeroWrite	Waiting to fill a dynamic shared memory backing file with zeroes.
LockFileAddtodatadirRead	Waiting for a read while adding a line to the data directory lock file.
LockFileAddtodatadirSync	Waiting for data to reach durable storage while adding a line to the data directory lock file.
LockFileAddtodatadirWrite	Waiting for a write while adding a line to the data directory lock file.
LockFileCreateRead	Waiting to read while creating the data directory lock file.
LockFileCreateSync	Waiting for data to reach durable storage while creating the data directory lock file.
LockFileCreateWrite	Waiting for a write while creating the data directory lock file.
LockFileRecheckdatadirRead	Waiting for a read during recheck of the data directory lock file.
LogicalRewriteCheckpointSync	Waiting for logical rewrite mappings to reach durable storage during a checkpoint.
LogicalRewriteMappingSync	Waiting for mapping data to reach durable storage during a logical rewrite.
LogicalRewriteMappingWrite	Waiting for a write of mapping data during a logical rewrite.

IO Wait Event	Description
LogicalRewriteSync	Waiting for logical rewrite mappings to reach durable storage.
LogicalRewriteTruncate	Waiting for truncate of mapping data during a logical rewrite.
LogicalRewriteWrite	Waiting for a write of logical rewrite mappings.
RelationMapRead	Waiting for a read of the relation map file.
RelationMapReplace	Waiting for durable replacement of a relation map file.
RelationMapWrite	Waiting for a write to the relation map file.
ReorderBufferRead	Waiting for a read during reorder buffer management.
ReorderBufferWrite	Waiting for a write during reorder buffer management.
ReorderLogicalMappingRead	Waiting for a read of a logical mapping during reorder buffer management.
ReplicationSlotRead	Waiting for a read from a replication slot control file.
ReplicationSlotRestoreSync	Waiting for a replication slot control file to reach durable storage while restoring it to memory.
ReplicationSlotSync	Waiting for a replication slot control file to reach durable storage.
ReplicationSlotWrite	Waiting for a write to a replication slot control file.
SlruFlushSync	Waiting for SLRU data to reach durable storage during a checkpoint or database shutdown.
SlruRead	Waiting for a read of an SLRU page.
SlruSync	Waiting for SLRU data to reach durable storage following a page write.
SlruWrite	Waiting for a write of an SLRU page.
SnapbuildRead	Waiting for a read of a serialized historical catalog snapshot.
SnapbuildSync	Waiting for a serialized historical catalog snapshot to reach durable storage.
SnapbuildWrite	Waiting for a write of a serialized historical catalog snapshot.
TimelineHistoryFileSync	Waiting for a timeline history file received via streaming replication to reach durable storage.
TimelineHistoryFileWrite	Waiting for a write of a timeline history file received via streaming replication.
TimelineHistoryRead	Waiting for a read of a timeline history file.
TimelineHistorySync	Waiting for a newly created timeline history file to reach durable storage.
TimelineHistoryWrite	Waiting for a write of a newly created timeline history file.
TwophaseFileRead	Waiting for a read of a two phase state file.

IO Wait Event	Description
TwophaseFileSync	Waiting for a two phase state file to reach durable storage.
TwophaseFileWrite	Waiting for a write of a two phase state file.
VersionFileSync	Waiting for the version file to reach durable storage while creating a database.
VersionFileWrite	Waiting for the version file to be written while creating a database.
WalsenderTimelineHistoryRead	Waiting for a read from a timeline history file during a walsender timeline command.
WalBootstrapSync	Waiting for WAL to reach durable storage during bootstrapping.
WalBootstrapWrite	Waiting for a write of a WAL page during bootstrapping.
WalCopyRead	Waiting for a read when creating a new WAL segment by copying an existing one.
WalCopySync	Waiting for a new WAL segment created by copying an existing one to reach durable storage.
WalCopyWrite	Waiting for a write when creating a new WAL segment by copying an existing one.
WalInitSync	Waiting for a newly initialized WAL file to reach durable storage.
WalInitWrite	Waiting for a write while initializing a new WAL file.
WalRead	Waiting for a read from a WAL file.
WalSummaryRead	Waiting for a read from a WAL summary file.
WalSummaryWrite	Waiting for a write to a WAL summary file.
WalSync	Waiting for a WAL file to reach durable storage.
WalSyncMethodAssign	Waiting for data to reach durable storage while assigning a new WAL sync method.
WalWrite	Waiting for a write to a WAL file.

Table 27.10. Wait Events of Type Ipc

IPC Wait Event	Description
AppendReady	Waiting for subplan nodes of an Append plan node to be ready.
ArchiveCleanupCommand	Waiting for archive_cleanup_command to complete.
ArchiveCommand	Waiting for archive_command to complete.
BackendTermination	Waiting for the termination of another backend.
BackupWaitWalArchive	Waiting for WAL files required for a backup to be successfully archived.
BgworkerShutdown	Waiting for background worker to shut down.
BgworkerStartup	Waiting for background worker to start up.
BtreePage	Waiting for the page number needed to continue a parallel B-tree scan to become available.
BufferIo	Waiting for buffer I/O to complete.

IPC Wait Event	Description
CheckpointDelayComplete	Waiting for a backend that blocks a checkpoint from completing.
CheckpointDelayStart	Waiting for a backend that blocks a checkpoint from starting.
CheckpointDone	Waiting for a checkpoint to complete.
CheckpointStart	Waiting for a checkpoint to start.
ExecuteGather	Waiting for activity from a child process while executing a <code>Gather</code> plan node.
HashBatchAllocate	Waiting for an elected Parallel Hash participant to allocate a hash table.
HashBatchElect	Waiting to elect a Parallel Hash participant to allocate a hash table.
HashBatchLoad	Waiting for other Parallel Hash participants to finish loading a hash table.
HashBuildAllocate	Waiting for an elected Parallel Hash participant to allocate the initial hash table.
HashBuildElect	Waiting to elect a Parallel Hash participant to allocate the initial hash table.
HashBuildHashInner	Waiting for other Parallel Hash participants to finish hashing the inner relation.
HashBuildHashOuter	Waiting for other Parallel Hash participants to finish partitioning the outer relation.
HashGrowBatchesDecide	Waiting to elect a Parallel Hash participant to decide on future batch growth.
HashGrowBatchesElect	Waiting to elect a Parallel Hash participant to allocate more batches.
HashGrowBatchesFinish	Waiting for an elected Parallel Hash participant to decide on future batch growth.
HashGrowBatchesReallocate	Waiting for an elected Parallel Hash participant to allocate more batches.
HashGrowBatchesRepartition	Waiting for other Parallel Hash participants to finish repartitioning.
HashGrowBucketsElect	Waiting to elect a Parallel Hash participant to allocate more buckets.
HashGrowBucketsReallocate	Waiting for an elected Parallel Hash participant to finish allocating more buckets.
HashGrowBucketsReinsert	Waiting for other Parallel Hash participants to finish inserting tuples into new buckets.
LogicalApplySendData	Waiting for a logical replication leader apply process to send data to a parallel apply process.
LogicalParallelApplyStateChange	Waiting for a logical replication parallel apply process to change state.
LogicalSyncData	Waiting for a logical replication remote server to send data for initial table synchronization.
LogicalSyncStateChange	Waiting for a logical replication remote server to change state.

IPC Wait Event	Description
MessageQueueInternal	Waiting for another process to be attached to a shared message queue.
MessageQueuePutMessage	Waiting to write a protocol message to a shared message queue.
MessageQueueReceive	Waiting to receive bytes from a shared message queue.
MessageQueueSend	Waiting to send bytes to a shared message queue.
MultixactCreation	Waiting for a multixact creation to complete.
ParallelBitmapScan	Waiting for parallel bitmap scan to become initialized.
ParallelCreateIndexScan	Waiting for parallel <code>CREATE INDEX</code> workers to finish heap scan.
ParallelFinish	Waiting for parallel workers to finish computing.
ProccarrayGroupUpdate	Waiting for the group leader to clear the transaction ID at transaction end.
ProcSignalBarrier	Waiting for a barrier event to be processed by all backends.
Promote	Waiting for standby promotion.
RecoveryConflictSnapshot	Waiting for recovery conflict resolution for a vacuum cleanup.
RecoveryConflictTablespace	Waiting for recovery conflict resolution for dropping a tablespace.
RecoveryEndCommand	Waiting for recovery_end_command to complete.
RecoveryPause	Waiting for recovery to be resumed.
ReplicationOriginDrop	Waiting for a replication origin to become inactive so it can be dropped.
ReplicationSlotDrop	Waiting for a replication slot to become inactive so it can be dropped.
RestoreCommand	Waiting for restore_command to complete.
SafeSnapshot	Waiting to obtain a valid snapshot for a <code>READ ONLY DEFERRABLE</code> transaction.
SyncRep	Waiting for confirmation from a remote server during synchronous replication.
WalReceiverExit	Waiting for the WAL receiver to exit.
WalReceiverWaitStart	Waiting for startup process to send initial data for streaming replication.
WalSummaryReady	Waiting for a new WAL summary to be generated.
XactGroupUpdate	Waiting for the group leader to update transaction status at transaction end.

Table 27.11. Wait Events of Type `Lock`

Lock Wait Event	Description
advisory	Waiting to acquire an advisory user lock.
applytransaction	Waiting to acquire a lock on a remote transaction being applied by a logical replication subscriber.
extend	Waiting to extend a relation.

Lock Wait Event	Description
frozenid	Waiting to update <code>pg_database.datfrozenxid</code> and <code>pg_database.datminmxid</code> .
object	Waiting to acquire a lock on a non-relation database object.
page	Waiting to acquire a lock on a page of a relation.
relation	Waiting to acquire a lock on a relation.
spectoken	Waiting to acquire a speculative insertion lock.
transactionid	Waiting for a transaction to finish.
tuple	Waiting to acquire a lock on a tuple.
userlock	Waiting to acquire a user lock.
virtualxid	Waiting to acquire a virtual transaction ID lock; see Section 67.1 .

Table 27.12. Wait Events of Type `Lwlock`

LWLock Wait Event	Description
AddinShmemInit	Waiting to manage an extension's space allocation in shared memory.
AutoFile	Waiting to update the <code>postgresql.auto.conf</code> file.
Autovacuum	Waiting to read or update the current state of autovacuum workers.
AutovacuumSchedule	Waiting to ensure that a table selected for autovacuum still needs vacuuming.
BackgroundWorker	Waiting to read or update background worker state.
BtreeVacuum	Waiting to read or update vacuum-related information for a B-tree index.
BufferContent	Waiting to access a data page in memory.
BufferMapping	Waiting to associate a data block with a buffer in the buffer pool.
CfsGc	Waiting on CFS GC control lock.
CheckpointInterComm	Waiting to manage fsync requests.
CommitTs	Waiting to read or update the last value set for a transaction commit timestamp.
CommitTsBuffer	Waiting for I/O on a commit timestamp SLRU buffer.
CommitTsSLRU	Waiting to access the commit timestamp SLRU cache.
ControlFile	Waiting to read or update the <code>pg_control</code> file or create a new WAL file.
CSNLogSLRU	Waiting to access CSN SLRU cache.
CSNSLogWal	Waiting to access CSN/LSN map.
CSNSnapshotXidMap	Waiting to access CSN Snapshot/Xid Map.
DSMRegistry	Waiting to read or update the dynamic shared memory registry.

LWLock Wait Event	Description
DSMRegistryDSA	Waiting to access dynamic shared memory registry's dynamic shared memory allocator.
DSMRegistryHash	Waiting to access dynamic shared memory registry's shared hash table.
DynamicSharedMemoryControl	Waiting to read or update dynamic shared memory allocation information.
InjectionPoint	Waiting to read or update information related to injection points.
LockFastPath	Waiting to read or update a process' fast-path lock information.
LockManager	Waiting to read or update information about "heavyweight" locks.
LogicalRepLauncherDSA	Waiting to access logical replication launcher's dynamic shared memory allocator.
LogicalRepLauncherHash	Waiting to access logical replication launcher's shared hash table.
LogicalRepWorker	Waiting to read or update the state of logical replication workers.
MultiXactGen	Waiting to read or update shared multixact state.
MultiXactMemberBuffer	Waiting for I/O on a multixact member SLRU buffer.
MultiXactMemberSLRU	Waiting to access the multixact member SLRU cache.
MultiXactOffsetBuffer	Waiting for I/O on a multixact offset SLRU buffer.
MultiXactOffsetSLRU	Waiting to access the multixact offset SLRU cache.
MultiXactTruncation	Waiting to read or truncate multixact information.
NotifyBuffer	Waiting for I/O on a NOTIFY message SLRU buffer.
NotifyQueue	Waiting to read or update NOTIFY messages.
NotifyQueueTail	Waiting to update limit on NOTIFY message storage.
NotifySLRU	Waiting to access the NOTIFY message SLRU cache.
OidGen	Waiting to allocate a new OID.
ParallelAppend	Waiting to choose the next subplan during Parallel Append plan execution.
ParallelHashJoin	Waiting to synchronize workers during Parallel Hash Join plan execution.
ParallelQueryDSA	Waiting for parallel query dynamic shared memory allocation.
ParallelVacuumDSA	Waiting for parallel vacuum dynamic shared memory allocation.
PerSessionDSA	Waiting for parallel query dynamic shared memory allocation.
PerSessionRecordType	Waiting to access a parallel query's information about composite types.

LWLock Wait Event	Description
PerSessionRecordTypmod	Waiting to access a parallel query's information about type modifiers that identify anonymous record types.
PerXactPredicateList	Waiting to access the list of predicate locks held by the current serializable transaction during a parallel query.
PgStatsData	Waiting for shared memory stats data access.
PgStatsDSA	Waiting for stats dynamic shared memory allocator access.
PgStatsHash	Waiting for stats shared memory hash table access.
PredicateLockManager	Waiting to access predicate lock information used by serializable transactions.
ProcArray	Waiting to access the shared per-process data structures (typically, to get a snapshot or report a session's transaction ID).
RelationMapping	Waiting to read or update a <code>pg_filenode.map</code> file (used to track the filenode assignments of certain system catalogs).
RelCacheInit	Waiting to read or update a <code>pg_internal.init</code> relation cache initialization file.
ReplicationOrigin	Waiting to create, drop or use a replication origin.
ReplicationOriginState	Waiting to read or update the progress of one replication origin.
ReplicationSlotAllocation	Waiting to allocate or free a replication slot.
ReplicationSlotControl	Waiting to read or update replication slot state.
ReplicationSlotIO	Waiting for I/O on a replication slot.
SerialBuffer	Waiting for I/O on a serializable transaction conflict SLRU buffer.
SerialControl	Waiting to read or update shared <code>pg_serial</code> state.
SerializableFinishedList	Waiting to access the list of finished serializable transactions.
SerializablePredicateList	Waiting to access the list of predicate locks held by serializable transactions.
SerializableXactHash	Waiting to read or update information about serializable transactions.
SerialSLRU	Waiting to access the serializable transaction conflict SLRU cache.
SharedTidBitmap	Waiting to access a shared TID bitmap during a parallel bitmap index scan.
SharedTupleStore	Waiting to access a shared tuple store during parallel query.
ShmemIndex	Waiting to find or allocate space in shared memory.
SInvalRead	Waiting to retrieve messages from the shared catalog invalidation queue.

LWLock Wait Event	Description
SInvalWrite	Waiting to add a message to the shared catalog invalidation queue.
SubtransBuffer	Waiting for I/O on a sub-transaction SLRU buffer.
SubtransSLRU	Waiting to access the sub-transaction SLRU cache.
SyncRep	Waiting to read or update information about the state of synchronous replication.
SyncScan	Waiting to select the starting location of a synchronized table scan.
TablespaceCreate	Waiting to create or drop a tablespace.
TwoPhaseState	Waiting to read or update the state of prepared transactions.
WaitEventCustom	Waiting to read or update custom wait events information.
WALBufMapping	Waiting to replace a page in WAL buffers.
WALInsert	Waiting to insert WAL data into a memory buffer.
WALSummarizer	Waiting to read or update WAL summarization state.
WALWrite	Waiting for WAL buffers to be written to disk.
WrapLimitsVacuum	Waiting to update limits on transaction id and multixact consumption.
XactBuffer	Waiting for I/O on a transaction status SLRU buffer.
XactSLRU	Waiting to access the transaction status SLRU cache.
XactTruncation	Waiting to execute <code>pg_xact_status</code> or update the oldest transaction ID available to it.
XidGen	Waiting to allocate a new transaction ID.

Table 27.13. Wait Events of Type Timeout

Timeout Wait Event	Description
BaseBackupThrottle	Waiting during base backup when throttling activity.
CheckpointWriteDelay	Waiting between writes while performing a checkpoint.
PgSleep	Waiting due to a call to <code>pg_sleep</code> or a sibling function.
RecoveryApplyDelay	Waiting to apply WAL during recovery because of a delay setting.
RecoveryRetrieveRetryInterval	Waiting during recovery when WAL data is not available from any source (<code>pg_wal</code> , archive or stream).
RegisterSyncRequest	Waiting while sending synchronization requests to the checkpointer, because the request queue is full.
SpinDelay	Waiting while acquiring a contended spinlock.

Timeout Wait Event	Description
VacuumDelay	Waiting in a cost-based vacuum delay point.
VacuumTruncate	Waiting to acquire an exclusive lock to truncate off any empty pages at the end of a table vacuumed.
WalSummarizerError	Waiting after a WAL summarizer error.

Here are examples of how wait events can be viewed:

```
SELECT pid, wait_event_type, wait_event FROM pg_stat_activity WHERE wait_event is NOT NULL;
```

```
pid | wait_event_type | wait_event
-----+-----+-----
2540 | Lock            | relation
6644 | LWLock          | ProcArray
(2 rows)
```

```
SELECT a.pid, a.wait_event, w.description
FROM pg_stat_activity a JOIN
     pg_wait_events w ON (a.wait_event_type = w.type AND
                        a.wait_event = w.name)
WHERE a.wait_event is NOT NULL and a.state = 'active';
```

```
-[ RECORD 1 ]-----
pid          | 686674
wait_event    | WALInitSync
description   | Waiting for a newly initialized WAL file to reach durable storage
```

Note

Extensions can add `Extension`, `InjectionPoint`, and `LWLock` events to the lists shown in [Table 27.8](#) and [Table 27.12](#). In some cases, the name of an `LWLock` assigned by an extension will not be available in all server processes. It might be reported as just “extension” rather than the extension-assigned name.

27.2.4. pgpro_stat_wal_activity

The `pgpro_stat_wal_activity` view will have one row per server process, showing information related to the current activity of that process. It inherits all columns from the [pg_stat_activity](#) view and adds the `wal_written` column to display the size of WAL files generated by each process.

Table 27.14. pgpro_stat_wal_activity View

Column	Type	Description
wal_written	bigint	The size of WAL files generated by a process.

27.2.5. pg_stat_replication

The `pg_stat_replication` view will contain one row per WAL sender process, showing statistics about replication to that sender's connected standby server. Only directly connected standbys are listed; no information is available about downstream standby servers.

Table 27.15. pg_stat_replication View

Column	Type	Description
pid	integer	Process ID of a WAL sender process

Column	Type	Description
usesysid	oid	OID of the user logged into this WAL sender process
username	name	Name of the user logged into this WAL sender process
application_name	text	Name of the application that is connected to this WAL sender
client_addr	inet	IP address of the client connected to this WAL sender. If this field is null, it indicates that the client is connected via a Unix socket on the server machine.
client_hostname	text	Host name of the connected client, as reported by a reverse DNS lookup of <code>client_addr</code> . This field will only be non-null for IP connections, and only when log_hostname is enabled.
client_port	integer	TCP port number that the client is using for communication with this WAL sender, or -1 if a Unix socket is used
backend_start	timestamp with time zone	Time when this process was started, i.e., when the client connected to this WAL sender
backend_xmin	xid	This standby's xmin horizon reported by hot_standby_feedback .
state	text	Current WAL sender state. Possible values are: <ul style="list-style-type: none"> <code>startup</code>: This WAL sender is starting up. <code>catchup</code>: This WAL sender's connected standby is catching up with the primary. <code>streaming</code>: This WAL sender is streaming changes after its connected standby server has caught up with the primary. <code>backup</code>: This WAL sender is sending a backup. <code>stopping</code>: This WAL sender is stopping.
sent_lsn	pg_lsn	Last write-ahead log location sent on this connection
write_lsn	pg_lsn	Last write-ahead log location written to disk by this standby server
flush_lsn	pg_lsn	Last write-ahead log location flushed to disk by this standby server
replay_lsn	pg_lsn	Last write-ahead log location replayed into the database on this standby server
write_lag	interval	Time elapsed between flushing recent WAL locally and receiving notification that this standby server has written it (but not yet flushed it or applied it). This can be used to gauge the delay that <code>synchronous_commit level remote_write</code> incurred while committing if this server was configured as a synchronous standby.
flush_lag	interval	Time elapsed between flushing recent WAL locally and receiving notification that this standby server has written and flushed it (but not yet applied it). This can be used to gauge the delay that <code>synchronous_commit level on</code> incurred while committing if this server was configured as a synchronous standby.
replay_lag	interval	

Column Type	Description
<code>nous_commit_level remote_apply</code>	Time elapsed between flushing recent WAL locally and receiving notification that this standby server has written, flushed and applied it. This can be used to gauge the delay that <code>synchronous_commit level remote_apply</code> incurred while committing if this server was configured as a synchronous standby.
<code>sync_priority integer</code>	Priority of this standby server for being chosen as the synchronous standby in a priority-based synchronous replication. This has no effect in a quorum-based synchronous replication.
<code>sync_state text</code>	Synchronous state of this standby server. Possible values are: <ul style="list-style-type: none"> <code>async</code>: This standby server is asynchronous. <code>potential</code>: This standby server is now asynchronous, but can potentially become synchronous if one of current synchronous ones fails. <code>sync</code>: This standby server is synchronous. <code>quorum</code>: This standby server is considered as a candidate for quorum standbys.
<code>reply_time timestamp with time zone</code>	Send time of last reply message received from standby server

The lag times reported in the `pg_stat_replication` view are measurements of the time taken for recent WAL to be written, flushed and replayed and for the sender to know about it. These times represent the commit delay that was (or would have been) introduced by each synchronous commit level, if the remote server was configured as a synchronous standby. For an asynchronous standby, the `replay_lag` column approximates the delay before recent transactions became visible to queries. If the standby server has entirely caught up with the sending server and there is no more WAL activity, the most recently measured lag times will continue to be displayed for a short time and then show NULL.

Lag times work automatically for physical replication. Logical decoding plugins may optionally emit tracking messages; if they do not, the tracking mechanism will simply display NULL lag.

Note

The reported lag times are not predictions of how long it will take for the standby to catch up with the sending server assuming the current rate of replay. Such a system would show similar times while new WAL is being generated, but would differ when the sender becomes idle. In particular, when the standby has caught up completely, `pg_stat_replication` shows the time taken to write, flush and replay the most recent reported WAL location rather than zero as some users might expect. This is consistent with the goal of measuring synchronous commit and transaction visibility delays for recent write transactions. To reduce confusion for users expecting a different model of lag, the lag columns revert to NULL after a short time on a fully replayed idle system. Monitoring systems should choose whether to represent this as missing data, zero or continue to display the last known value.

27.2.6. `pg_stat_replication_slots`

The `pg_stat_replication_slots` view will contain one row per logical replication slot, showing statistics about its usage.

Table 27.16. `pg_stat_replication_slots` View

Column Type	Description
<code>slot_name text</code>	

Column Type	Description
	A unique, cluster-wide identifier for the replication slot
<code>spill_txns bigint</code>	Number of transactions spilled to disk once the memory used by logical decoding to decode changes from WAL has exceeded <code>logical_decoding_work_mem</code> . The counter gets incremented for both top-level transactions and subtransactions.
<code>spill_count bigint</code>	Number of times transactions were spilled to disk while decoding changes from WAL for this slot. This counter is incremented each time a transaction is spilled, and the same transaction may be spilled multiple times.
<code>spill_bytes bigint</code>	Amount of decoded transaction data spilled to disk while performing decoding of changes from WAL for this slot. This and other spill counters can be used to gauge the I/O which occurred during logical decoding and allow tuning <code>logical_decoding_work_mem</code> .
<code>stream_txns bigint</code>	Number of in-progress transactions streamed to the decoding output plugin after the memory used by logical decoding to decode changes from WAL for this slot has exceeded <code>logical_decoding_work_mem</code> . Streaming only works with top-level transactions (subtransactions can't be streamed independently), so the counter is not incremented for subtransactions.
<code>stream_count bigint</code>	Number of times in-progress transactions were streamed to the decoding output plugin while decoding changes from WAL for this slot. This counter is incremented each time a transaction is streamed, and the same transaction may be streamed multiple times.
<code>stream_bytes bigint</code>	Amount of transaction data decoded for streaming in-progress transactions to the decoding output plugin while decoding changes from WAL for this slot. This and other streaming counters for this slot can be used to tune <code>logical_decoding_work_mem</code> .
<code>total_txns bigint</code>	Number of decoded transactions sent to the decoding output plugin for this slot. This counts top-level transactions only, and is not incremented for subtransactions. Note that this includes the transactions that are streamed and/or spilled.
<code>total_bytes bigint</code>	Amount of transaction data decoded for sending transactions to the decoding output plugin while decoding changes from WAL for this slot. Note that this includes data that is streamed and/or spilled.
<code>stats_reset timestamp with time zone</code>	Time at which these statistics were last reset

27.2.7. `pg_stat_wal_receiver`

The `pg_stat_wal_receiver` view will contain only one row, showing statistics about the WAL receiver from that receiver's connected server.

Table 27.17. `pg_stat_wal_receiver` View

Column Type	Description
<code>pid integer</code>	Process ID of the WAL receiver process
<code>status text</code>	Activity status of the WAL receiver process
<code>receive_start_lsn pg_lsn</code>	

Column Type	Description
	First write-ahead log location used when WAL receiver is started
<code>receive_start_tli</code> integer	First timeline number used when WAL receiver is started
<code>written_lsn</code> <code>pg_lsn</code>	Last write-ahead log location already received and written to disk, but not flushed. This should not be used for data integrity checks.
<code>flushed_lsn</code> <code>pg_lsn</code>	Last write-ahead log location already received and flushed to disk, the initial value of this field being the first log location used when WAL receiver is started
<code>received_tli</code> integer	Timeline number of last write-ahead log location received and flushed to disk, the initial value of this field being the timeline number of the first log location used when WAL receiver is started
<code>last_msg_send_time</code> timestamp with time zone	Send time of last message received from origin WAL sender
<code>last_msg_receipt_time</code> timestamp with time zone	Receipt time of last message received from origin WAL sender
<code>latest_end_lsn</code> <code>pg_lsn</code>	Last write-ahead log location reported to origin WAL sender
<code>latest_end_time</code> timestamp with time zone	Time of last write-ahead log location reported to origin WAL sender
<code>slot_name</code> text	Replication slot name used by this WAL receiver
<code>sender_host</code> text	Host of the Postgres Pro Shardman instance this WAL receiver is connected to. This can be a host name, an IP address, or a directory path if the connection is via Unix socket. (The path case can be distinguished because it will always be an absolute path, beginning with /.)
<code>sender_port</code> integer	Port number of the Postgres Pro Shardman instance this WAL receiver is connected to.
<code>conninfo</code> text	Connection string used by this WAL receiver, with security-sensitive fields obfuscated.

27.2.8. `pg_stat_recovery_prefetch`

The `pg_stat_recovery_prefetch` view will contain only one row. The columns `wal_distance`, `block_distance` and `io_depth` show current values, and the other columns show cumulative counters that can be reset with the `pg_stat_reset_shared` function.

Table 27.18. `pg_stat_recovery_prefetch` View

Column Type	Description
<code>stats_reset</code> timestamp with time zone	Time at which these statistics were last reset
<code>prefetch</code> bigint	Number of blocks prefetched because they were not in the buffer pool
<code>hit</code> bigint	Number of blocks not prefetched because they were already in the buffer pool
<code>skip_init</code> bigint	

Column	Type	Description
		Number of blocks not prefetched because they would be zero-initialized
skip_new	bigint	Number of blocks not prefetched because they didn't exist yet
skip_fpw	bigint	Number of blocks not prefetched because a full page image was included in the WAL
skip_rep	bigint	Number of blocks not prefetched because they were already recently prefetched
wal_distance	int	How many bytes ahead the prefetcher is looking
block_distance	int	How many blocks ahead the prefetcher is looking
io_depth	int	How many prefetches have been initiated but are not yet known to have completed

27.2.9. pg_stat_subscription

Table 27.19. pg_stat_subscription View

Column	Type	Description
subid	oid	OID of the subscription
subname	name	Name of the subscription
worker_type	text	Type of the subscription worker process. Possible types are <code>apply</code> , <code>parallel apply</code> , and <code>table synchronization</code> .
pid	integer	Process ID of the subscription worker process
leader_pid	integer	Process ID of the leader apply worker if this process is a parallel apply worker; NULL if this process is a leader apply worker or a table synchronization worker
relid	oid	OID of the relation that the worker is synchronizing; NULL for the leader apply worker and parallel apply workers
received_lsn	pg_lsn	Last write-ahead log location received, the initial value of this field being 0; NULL for parallel apply workers
last_msg_send_time	timestamp with time zone	Send time of last message received from origin WAL sender; NULL for parallel apply workers
last_msg_receipt_time	timestamp with time zone	Receipt time of last message received from origin WAL sender; NULL for parallel apply workers
latest_end_lsn	pg_lsn	Last write-ahead log location reported to origin WAL sender; NULL for parallel apply workers
latest_end_time	timestamp with time zone	Time of last write-ahead log location reported to origin WAL sender; NULL for parallel apply workers

27.2.10. pg_stat_subscription_stats

The `pg_stat_subscription_stats` view will contain one row per subscription.

Table 27.20. pg_stat_subscription_stats View

Column	Type	Description
subid	oid	OID of the subscription
subname	name	Name of the subscription
apply_error_count	bigint	Number of times an error occurred while applying changes
sync_error_count	bigint	Number of times an error occurred during the initial table synchronization
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.11. pg_stat_ssl

The `pg_stat_ssl` view will contain one row per backend or WAL sender process, showing statistics about SSL usage on this connection. It can be joined to `pg_stat_activity` or `pg_stat_replication` on the `pid` column to get more details about the connection.

Table 27.21. pg_stat_ssl View

Column	Type	Description
pid	integer	Process ID of a backend or WAL sender process
ssl	boolean	True if SSL is used on this connection
version	text	Version of SSL in use, or NULL if SSL is not in use on this connection
cipher	text	Name of SSL cipher in use, or NULL if SSL is not in use on this connection
bits	integer	Number of bits in the encryption algorithm used, or NULL if SSL is not used on this connection
client_dn	text	Distinguished Name (DN) field from the client certificate used, or NULL if no client certificate was supplied or if SSL is not in use on this connection. This field is truncated if the DN field is longer than NAMEDATALEN (64 characters in a standard build).
client_serial	numeric	Serial number of the client certificate, or NULL if no client certificate was supplied or if SSL is not in use on this connection. The combination of certificate serial number and certificate issuer uniquely identifies a certificate (unless the issuer erroneously reuses serial numbers).
issuer_dn	text	DN of the issuer of the client certificate, or NULL if no client certificate was supplied or if SSL is not in use on this connection. This field is truncated like <code>client_dn</code> .

27.2.12. pg_stat_gssapi

The `pg_stat_gssapi` view will contain one row per backend, showing information about GSSAPI usage on this connection. It can be joined to `pg_stat_activity` or `pg_stat_replication` on the `pid` column to get more details about the connection.

Table 27.22. `pg_stat_gssapi` View

Column	Type	Description
<code>pid</code>	integer	Process ID of a backend
<code>gss_authenticated</code>	boolean	True if GSSAPI authentication was used for this connection
<code>principal</code>	text	Principal used to authenticate this connection, or NULL if GSSAPI was not used to authenticate this connection. This field is truncated if the principal is longer than <code>NAMEDATALEN</code> (64 characters in a standard build).
<code>encrypted</code>	boolean	True if GSSAPI encryption is in use on this connection
<code>credentials_delegated</code>	boolean	True if GSSAPI credentials were delegated on this connection.

27.2.13. `pg_stat_archiver`

The `pg_stat_archiver` view will always have a single row, containing data about the archiver process of the cluster.

Table 27.23. `pg_stat_archiver` View

Column	Type	Description
<code>archived_count</code>	bigint	Number of WAL files that have been successfully archived
<code>last_archived_wal</code>	text	Name of the WAL file most recently successfully archived
<code>last_archived_time</code>	timestamp with time zone	Time of the most recent successful archive operation
<code>failed_count</code>	bigint	Number of failed attempts for archiving WAL files
<code>last_failed_wal</code>	text	Name of the WAL file of the most recent failed archival operation
<code>last_failed_time</code>	timestamp with time zone	Time of the most recent failed archival operation
<code>stats_reset</code>	timestamp with time zone	Time at which these statistics were last reset

Normally, WAL files are archived in order, oldest to newest, but that is not guaranteed, and does not hold under special circumstances like when promoting a standby or after crash recovery. Therefore it is not safe to assume that all files older than `last_archived_wal` have also been successfully archived.

27.2.14. `pg_stat_io`

The `pg_stat_io` view will contain one row for each combination of backend type, target I/O object, and I/O context, showing cluster-wide I/O statistics. Combinations which do not make sense are omitted.

Currently, I/O on relations (e.g. tables, indexes) is tracked. However, relation I/O which bypasses shared buffers (e.g. when moving a table from one tablespace to another) is currently not tracked.

Table 27.24. `pg_stat_io` View

Column	Type	Description
<code>backend_type</code>	text	Type of backend (e.g. background worker, autovacuum worker). See pg_stat_activity for more information on <code>backend_type</code> s. Some <code>backend_type</code> s do not accumulate I/O operation statistics and will not be included in the view.
<code>object</code>	text	Target object of an I/O operation. Possible values are: <ul style="list-style-type: none"> <code>relation</code>: Permanent relations. <code>temp relation</code>: Temporary relations.
<code>context</code>	text	The context of an I/O operation. Possible values are: <ul style="list-style-type: none"> <code>normal</code>: The default or standard <code>context</code> for a type of I/O operation. For example, by default, relation data is read into and written out from shared buffers. Thus, reads and writes of relation data to and from shared buffers are tracked in <code>context normal</code>. <code>vacuum</code>: I/O operations performed outside of shared buffers while vacuuming and analyzing permanent relations. Temporary table vacuums use the same local buffer pool as other temporary table I/O operations and are tracked in <code>context normal</code>. <code>bulkread</code>: Certain large read I/O operations done outside of shared buffers, for example, a sequential scan of a large table. <code>bulkwrite</code>: Certain large write I/O operations done outside of shared buffers, such as <code>COPY</code>.
<code>reads</code>	bigint	Number of read operations, each of the size specified in <code>op_bytes</code> .
<code>read_time</code>	double precision	Time spent in read operations in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>writes</code>	bigint	Number of write operations, each of the size specified in <code>op_bytes</code> .
<code>write_time</code>	double precision	Time spent in write operations in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>writebacks</code>	bigint	Number of units of size <code>op_bytes</code> which the process requested the kernel write out to permanent storage.
<code>writeback_time</code>	double precision	Time spent in writeback operations in milliseconds (if track_io_timing is enabled, otherwise zero). This includes the time spent queueing write-out requests and, potentially, the time spent to write out the dirty data.
<code>extends</code>	bigint	Number of relation extend operations, each of the size specified in <code>op_bytes</code> .
<code>extend_time</code>	double precision	Time spent in extend operations in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>op_bytes</code>	bigint	The number of bytes per unit of I/O read, written, or extended. Relation data reads, writes, and extends are done in <code>block_size</code> units, derived from the build-time parameter <code>BLCKSZ</code> , which is 8192 by default.
<code>hits</code>	bigint	

Column Type	Description
	The number of times a desired block was found in a shared buffer.
<code>evictions bigint</code>	Number of times a block has been written out from a shared or local buffer in order to make it available for another use. In <code>context normal</code> , this counts the number of times a block was evicted from a buffer and replaced with another block. In <code>contexts bulkwrite</code> , <code>bulkread</code> , and <code>vacuum</code> , this counts the number of times a block was evicted from shared buffers in order to add the shared buffer to a separate, size-limited ring buffer for use in a bulk I/O operation.
<code>reuses bigint</code>	The number of times an existing buffer in a size-limited ring buffer outside of shared buffers was reused as part of an I/O operation in the <code>bulkread</code> , <code>bulkwrite</code> , or <code>vacuum contexts</code> .
<code>fsyncs bigint</code>	Number of <code>fsync</code> calls. These are only tracked in <code>context normal</code> .
<code>fsync_time double precision</code>	Time spent in <code>fsync</code> operations in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>stats_reset timestamp with time zone</code>	Time at which these statistics were last reset.

Some backend types never perform I/O operations on some I/O objects and/or in some I/O contexts. These rows are omitted from the view. For example, the checkpointer does not checkpoint temporary tables, so there will be no rows for `backend_type checkpointer` and `object temp relation`.

In addition, some I/O operations will never be performed either by certain backend types or on certain I/O objects and/or in certain I/O contexts. These cells will be NULL. For example, temporary tables are not `fsynced`, so `fsyncs` will be NULL for `object temp relation`. Also, the background writer does not perform reads, so `reads` will be NULL in rows for `backend_type background writer`.

`pg_stat_io` can be used to inform database tuning. For example:

- A high `evictions` count can indicate that shared buffers should be increased.
- Client backends rely on the checkpointer to ensure data is persisted to permanent storage. Large numbers of `fsyncs` by `client backends` could indicate a misconfiguration of shared buffers or of the checkpointer. More information on configuring the checkpointer can be found in [Section 28.5](#).
- Normally, client backends should be able to rely on auxiliary processes like the checkpointer and the background writer to write out dirty data as much as possible. Large numbers of writes by client backends could indicate a misconfiguration of shared buffers or of the checkpointer. More information on configuring the checkpointer can be found in [Section 28.5](#).

Note

Columns tracking I/O time will only be non-zero when [track_io_timing](#) is enabled. The user should be careful when referencing these columns in combination with their corresponding I/O operations in case `track_io_timing` was not enabled for the entire time since the last stats reset.

27.2.15. pg_stat_bgwriter

The `pg_stat_bgwriter` view will always have a single row, containing data about the background writer of the cluster.

Table 27.25. `pg_stat_bgwriter` View

Column Type	Description
<code>buffers_clean bigint</code>	

Column Type	Description
	Number of buffers written by the background writer
maxwritten_clean bigint	Number of times the background writer stopped a cleaning scan because it had written too many buffers
buffers_alloc bigint	Number of buffers allocated
stats_reset timestamp with time zone	Time at which these statistics were last reset

27.2.16. pg_stat_checkpoint

The `pg_stat_checkpoint` view will always have a single row, containing data about the checkpoint process of the cluster.

Table 27.26. pg_stat_checkpoint View

Column Type	Description
num_timed bigint	Number of scheduled checkpoints due to timeout. Note that checkpoints may be skipped if the server has been idle since the last one, and this value counts both completed and skipped checkpoints
num_requested bigint	Number of requested checkpoints that have been performed
restartpoints_timed bigint	Number of scheduled restartpoints due to timeout or after a failed attempt to perform it
restartpoints_req bigint	Number of requested restartpoints
restartpoints_done bigint	Number of restartpoints that have been performed
write_time double precision	Total amount of time that has been spent in the portion of processing checkpoints and restartpoints where files are written to disk, in milliseconds
sync_time double precision	Total amount of time that has been spent in the portion of processing checkpoints and restartpoints where files are synchronized to disk, in milliseconds
buffers_written bigint	Number of buffers written during checkpoints and restartpoints
stats_reset timestamp with time zone	Time at which these statistics were last reset

27.2.17. pg_stat_wal

The `pg_stat_wal` view will always have a single row, containing data about WAL activity of the cluster.

Table 27.27. pg_stat_wal View

Column Type	Description
wal_records bigint	Total number of WAL records generated

Column Type	Description
wal_fpi bigint	Total number of WAL full page images generated
wal_bytes numeric	Total amount of WAL generated in bytes
wal_buffers_full bigint	Number of times WAL data was written to disk because WAL buffers became full
wal_write bigint	Number of times WAL buffers were written out to disk via XLogWrite request. See Section 28.5 for more information about the internal WAL function XLogWrite.
wal_sync bigint	Number of times WAL files were synced to disk via <code>issue_xlog_fsync</code> request (if <code>fsync</code> is on and <code>wal_sync_method</code> is either <code>fdatasync</code> , <code>fsync</code> or <code>fsync_writethrough</code> , otherwise zero). See Section 28.5 for more information about the internal WAL function <code>issue_xlog_fsync</code> .
wal_write_time double precision	Total amount of time spent writing WAL buffers to disk via XLogWrite request, in milliseconds (if <code>track_wal_io_timing</code> is enabled, otherwise zero). This includes the sync time when <code>wal_sync_method</code> is either <code>open_datasync</code> or <code>open_sync</code> .
wal_sync_time double precision	Total amount of time spent syncing WAL files to disk via <code>issue_xlog_fsync</code> request, in milliseconds (if <code>track_wal_io_timing</code> is enabled, <code>fsync</code> is on, and <code>wal_sync_method</code> is either <code>fdatasync</code> , <code>fsync</code> or <code>fsync_writethrough</code> , otherwise zero).
stats_reset timestamp with time zone	Time at which these statistics were last reset

27.2.18. pg_stat_database

The `pg_stat_database` view will contain one row for each database in the cluster, plus one for shared objects, showing database-wide statistics.

Table 27.28. pg_stat_database View

Column Type	Description
datid oid	OID of this database, or 0 for objects belonging to a shared relation
datname name	Name of this database, or NULL for shared objects.
numbackends integer	Number of backends currently connected to this database, or NULL for shared objects. This is the only column in this view that returns a value reflecting current state; all other columns return the accumulated values since the last reset.
xact_commit bigint	Number of transactions in this database that have been committed
xact_rollback bigint	Number of transactions in this database that have been rolled back
blks_read bigint	Number of disk blocks read in this database
blks_hit bigint	

Column Type	Description
	Number of times disk blocks were found already in the buffer cache, so that a read was not necessary (this only includes hits in the Postgres Pro buffer cache, not the operating system's file system cache)
tup_returned bigint	Number of live rows fetched by sequential scans and index entries returned by index scans in this database
tup_fetched bigint	Number of live rows fetched by index scans in this database
tup_inserted bigint	Number of rows inserted by queries in this database
tup_updated bigint	Number of rows updated by queries in this database
tup_deleted bigint	Number of rows deleted by queries in this database
conflicts bigint	Number of queries canceled due to conflicts with recovery in this database. (Conflicts occur only on standby servers; see pg_stat_database_conflicts for details.)
temp_files bigint	Number of temporary files created by queries in this database. All temporary files are counted, regardless of why the temporary file was created (e.g., sorting or hashing), and regardless of the log_temp_files setting.
temp_bytes bigint	Total amount of data written to temporary files by queries in this database. All temporary files are counted, regardless of why the temporary file was created, and regardless of the log_temp_files setting.
deadlocks bigint	Number of deadlocks detected in this database
checksum_failures bigint	Number of data page checksum failures detected in this database (or on a shared object), or NULL if data checksums are not enabled.
checksum_last_failure timestamp with time zone	Time at which the last data page checksum failure was detected in this database (or on a shared object), or NULL if data checksums are not enabled.
blk_read_time double precision	Time spent reading data file blocks by backends in this database, in milliseconds (if track_io_timing is enabled, otherwise zero)
blk_write_time double precision	Time spent writing data file blocks by backends in this database, in milliseconds (if track_io_timing is enabled, otherwise zero)
session_time double precision	Time spent by database sessions in this database, in milliseconds (note that statistics are only updated when the state of a session changes, so if sessions have been idle for a long time, this idle time won't be included)
active_time double precision	Time spent executing SQL statements in this database, in milliseconds (this corresponds to the states <code>active</code> and <code>fastpath</code> function call in pg_stat_activity)
idle_in_transaction_time double precision	

Column	Type	Description
		Time spent idling while in a transaction in this database, in milliseconds (this corresponds to the states <code>idle in transaction</code> and <code>idle in transaction (aborted)</code> in pg_stat_activity)
<code>sessions</code>	<code>bigint</code>	Total number of sessions established to this database
<code>sessions_abandoned</code>	<code>bigint</code>	Number of database sessions to this database that were terminated because connection to the client was lost
<code>sessions_fatal</code>	<code>bigint</code>	Number of database sessions to this database that were terminated by fatal errors
<code>sessions_killed</code>	<code>bigint</code>	Number of database sessions to this database that were terminated by operator intervention
<code>stats_reset</code>	<code>timestamp with time zone</code>	Time at which these statistics were last reset

27.2.19. `pg_stat_database_conflicts`

The `pg_stat_database_conflicts` view will contain one row per database, showing database-wide statistics about query cancels occurring due to conflicts with recovery on standby servers. This view will only contain information on standby servers, since conflicts do not occur on primary servers.

Table 27.29. `pg_stat_database_conflicts` View

Column	Type	Description
<code>datid</code>	<code>oid</code>	OID of a database
<code>datname</code>	<code>name</code>	Name of this database
<code>confl_tablespace</code>	<code>bigint</code>	Number of queries in this database that have been canceled due to dropped tablespaces
<code>confl_lock</code>	<code>bigint</code>	Number of queries in this database that have been canceled due to lock timeouts
<code>confl_snapshot</code>	<code>bigint</code>	Number of queries in this database that have been canceled due to old snapshots
<code>confl_bufferpin</code>	<code>bigint</code>	Number of queries in this database that have been canceled due to pinned buffers
<code>confl_deadlock</code>	<code>bigint</code>	Number of queries in this database that have been canceled due to deadlocks
<code>confl_active_logicalslot</code>	<code>bigint</code>	Number of uses of logical slots in this database that have been canceled due to old snapshots or too low a wal_level on the primary

27.2.20. `pg_stat_all_tables`

The `pg_stat_all_tables` view will contain one row for each table in the current database (including TOAST tables), showing statistics about accesses to that specific table. The `pg_stat_user_tables` and `pg_stat_sys_tables` views contain the same information, but filtered to only show user and system tables respectively.

Table 27.30. pg_stat_all_tables View

Column	Type	Description
relid	oid	OID of a table
schemaname	name	Name of the schema that this table is in
relname	name	Name of this table
seq_scan	bigint	Number of sequential scans initiated on this table
last_seq_scan	timestamp with time zone	The time of the last sequential scan on this table, based on the most recent transaction stop time
seq_tup_read	bigint	Number of live rows fetched by sequential scans
idx_scan	bigint	Number of index scans initiated on this table
last_idx_scan	timestamp with time zone	The time of the last index scan on this table, based on the most recent transaction stop time
idx_tup_fetch	bigint	Number of live rows fetched by index scans
n_tup_ins	bigint	Total number of rows inserted
n_tup_upd	bigint	Total number of rows updated. (This includes row updates counted in <code>n_tup_hot_upd</code> and <code>n_tup_newpage_upd</code> , and remaining non-HOT updates.)
n_tup_del	bigint	Total number of rows deleted
n_tup_hot_upd	bigint	Number of rows HOT updated . These are updates where no successor versions are required in indexes.
n_tup_newpage_upd	bigint	Number of rows updated where the successor version goes onto a <i>new</i> heap page, leaving behind an original version with a <code>t_ctid</code> field that points to a different heap page. These are always non-HOT updates.
n_live_tup	bigint	Estimated number of live rows
n_dead_tup	bigint	Estimated number of dead rows
n_mod_since_analyze	bigint	Estimated number of rows modified since this table was last analyzed
n_ins_since_vacuum	bigint	Estimated number of rows inserted since this table was last vacuumed
last_vacuum	timestamp with time zone	Last time at which this table was manually vacuumed (not counting <code>VACUUM FULL</code>)
last_autovacuum	timestamp with time zone	Last time at which this table was vacuumed by the autovacuum daemon

Column Type	Description
<code>last_analyze timestamp with time zone</code>	Last time at which this table was manually analyzed
<code>last_autoanalyze timestamp with time zone</code>	Last time at which this table was analyzed by the autovacuum daemon
<code>vacuum_count bigint</code>	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)
<code>autovacuum_count bigint</code>	Number of times this table has been vacuumed by the autovacuum daemon
<code>analyze_count bigint</code>	Number of times this table has been manually analyzed
<code>autoanalyze_count bigint</code>	Number of times this table has been analyzed by the autovacuum daemon

27.2.21. `pg_stat_all_indexes`

The `pg_stat_all_indexes` view will contain one row for each index in the current database, showing statistics about accesses to that specific index. The `pg_stat_user_indexes` and `pg_stat_sys_indexes` views contain the same information, but filtered to only show user and system indexes respectively.

Table 27.31. `pg_stat_all_indexes` View

Column Type	Description
<code>relid oid</code>	OID of the table for this index
<code>indexrelid oid</code>	OID of this index
<code>schemaname name</code>	Name of the schema this index is in
<code>relname name</code>	Name of the table for this index
<code>indexrelname name</code>	Name of this index
<code>idx_scan bigint</code>	Number of index scans initiated on this index
<code>last_idx_scan timestamp with time zone</code>	The time of the last scan on this index, based on the most recent transaction stop time
<code>idx_tup_read bigint</code>	Number of index entries returned by scans on this index
<code>idx_tup_fetch bigint</code>	Number of live table rows fetched by simple index scans using this index

Indexes can be used by simple index scans, “bitmap” index scans, and the optimizer. In a bitmap scan the output of several indexes can be combined via AND or OR rules, so it is difficult to associate individual heap row fetches with specific indexes when a bitmap scan is used. Therefore, a bitmap scan increments the `pg_stat_all_indexes.idx_tup_read` count(s) for the index(es) it uses, and it increments the `pg_stat_all_tables.idx_tup_fetch` count for the table, but it does not affect `pg_stat_all_indexes.idx_tup_fetch`. The optimizer also accesses indexes to check for supplied constants whose values are outside the recorded range of the optimizer statistics because the optimizer statistics might be stale.

Note

The `idx_tup_read` and `idx_tup_fetch` counts can be different even without any use of bitmap scans, because `idx_tup_read` counts index entries retrieved from the index while `idx_tup_fetch` counts live rows fetched from the table. The latter will be less if any dead or not-yet-committed rows are fetched using the index, or if any heap fetches are avoided by means of an index-only scan.

Note

Queries that use certain SQL constructs to search for rows matching any value out of a list or array of multiple scalar values (see [Section 9.25](#)) perform multiple “primitive” index scans (up to one primitive scan per scalar value) during query execution. Each internal primitive index scan increments `pg_stat_all_indexes.idx_scan`, so it's possible for the count of index scans to significantly exceed the total number of index scan executor node executions.

27.2.22. pg_statio_all_tables

The `pg_statio_all_tables` view will contain one row for each table in the current database (including TOAST tables), showing statistics about I/O on that specific table. The `pg_statio_user_tables` and `pg_statio_sys_tables` views contain the same information, but filtered to only show user and system tables respectively.

Table 27.32. pg_statio_all_tables View

Column	Type	Description
<code>relid</code>	<code>oid</code>	OID of a table
<code>schemaname</code>	<code>name</code>	Name of the schema that this table is in
<code>relname</code>	<code>name</code>	Name of this table
<code>heap_blks_read</code>	<code>bigint</code>	Number of disk blocks read from this table
<code>heap_blks_hit</code>	<code>bigint</code>	Number of buffer hits in this table
<code>idx_blks_read</code>	<code>bigint</code>	Number of disk blocks read from all indexes on this table
<code>idx_blks_hit</code>	<code>bigint</code>	Number of buffer hits in all indexes on this table
<code>toast_blks_read</code>	<code>bigint</code>	Number of disk blocks read from this table's TOAST table (if any)
<code>toast_blks_hit</code>	<code>bigint</code>	Number of buffer hits in this table's TOAST table (if any)
<code>tidx_blks_read</code>	<code>bigint</code>	Number of disk blocks read from this table's TOAST table indexes (if any)
<code>tidx_blks_hit</code>	<code>bigint</code>	Number of buffer hits in this table's TOAST table indexes (if any)

27.2.23. pg_statio_all_indexes

The `pg_statio_all_indexes` view will contain one row for each index in the current database, showing statistics about I/O on that specific index. The `pg_statio_user_indexes` and `pg_statio_sys_indexes` views contain the same information, but filtered to only show user and system indexes respectively.

Table 27.33. `pg_statio_all_indexes` View

Column Type	Description
<code>relid oid</code>	OID of the table for this index
<code>indexrelid oid</code>	OID of this index
<code>schemaname name</code>	Name of the schema this index is in
<code>relname name</code>	Name of the table for this index
<code>indexrelname name</code>	Name of this index
<code>idx_blks_read bigint</code>	Number of disk blocks read from this index
<code>idx_blks_hit bigint</code>	Number of buffer hits in this index

27.2.24. `pg_statio_all_sequences`

The `pg_statio_all_sequences` view will contain one row for each sequence in the current database, showing statistics about I/O on that specific sequence.

Table 27.34. `pg_statio_all_sequences` View

Column Type	Description
<code>relid oid</code>	OID of a sequence
<code>schemaname name</code>	Name of the schema this sequence is in
<code>relname name</code>	Name of this sequence
<code>blks_read bigint</code>	Number of disk blocks read from this sequence
<code>blks_hit bigint</code>	Number of buffer hits in this sequence

27.2.25. `pg_stat_user_functions`

The `pg_stat_user_functions` view will contain one row for each tracked function, showing statistics about executions of that function. The [track_functions](#) parameter controls exactly which functions are tracked.

Table 27.35. `pg_stat_user_functions` View

Column Type	Description
<code>funcid oid</code>	

Column	Type	Description
		OID of a function
schemaname	name	Name of the schema this function is in
funcname	name	Name of this function
calls	bigint	Number of times this function has been called
total_time	double precision	Total time spent in this function and all other functions called by it, in milliseconds
self_time	double precision	Total time spent in this function itself, not including other functions called by it, in milliseconds

27.2.26. pg_stat_slru

Postgres Pro Shardman accesses certain on-disk information via SLRU (*simple least-recently-used*) caches. The `pg_stat_slru` view will contain one row for each tracked SLRU cache, showing statistics about access to cached pages.

For each SLRU cache that's part of the core server, there is a configuration parameter that controls its size, with the suffix `_buffers` appended.

Table 27.36. `pg_stat_slru` View

Column	Type	Description
name	text	Name of the SLRU
blks_zeroed	bigint	Number of blocks zeroed during initializations
blks_hit	bigint	Number of times disk blocks were found already in the SLRU, so that a read was not necessary (this only includes hits in the SLRU, not the operating system's file system cache)
blks_read	bigint	Number of disk blocks read for this SLRU
blks_written	bigint	Number of disk blocks written for this SLRU
blks_exists	bigint	Number of blocks checked for existence for this SLRU
flushes	bigint	Number of flushes of dirty data for this SLRU
truncates	bigint	Number of truncates for this SLRU
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.27. Views for Distributed System

Postgres Pro Shardman provides a set of views to show statistics related to a distributed system's specifics.

27.2.27.1. shardman.pg_stat_csn

The `shardman.pg_stat_csn` view has one row showing statistics on delays that take place during import of CSN snapshots. These delays occur because system clocks on Postgres Pro Shardman cluster nodes may be out of sync. The delays negatively impact the performance by increasing the query latency. The `shardman.pg_stat_csn` view allows tracking these delays. The view data is based on the [cumulative statistics system](#). The columns of the view are shown in [Table 27.37](#).

Table 27.37. shardman.pg_stat_csn Columns

Name	Type	Description
<code>csn_snapshots_imported</code>	<code>bigint</code>	Total number of imported CSN snapshots
<code>csn_total_import_delay</code>	<code>interval</code>	Total duration of all delays in importing CSN snapshots, in microseconds
<code>csn_max_shift</code>	<code>bigint</code>	Maximum registered snapshot CSN shift that caused a delay
<code>local_oldest_csn</code>	<code>bigint</code>	CSN of the oldest transaction on the current node
<code>local_oldest_xid</code>	<code>xid</code>	XID of the oldest transaction on the current node
<code>indoubt_threshold_incidents</code>	<code>bigint</code>	Total number of transactions that exceeded the 10 seconds limit in the <code>inDoubt</code> state.
<code>CSNXidMap_head_csn</code>	<code>bigint</code>	Most recent CSN in the <code>CSNSnapshotXidMap</code>
<code>CSNXidMap_head_xid</code>	<code>xid</code>	XID corresponding to the most recent CSN in the <code>CSNSnapshotXidMap</code>
<code>CSNXidMap_tail_csn</code>	<code>bigint</code>	Oldest CSN in the <code>CSNSnapshotXidMap</code>
<code>CSNXidMap_tail_xid</code>	<code>xid</code>	XID corresponding to the oldest CSN in the <code>CSNSnapshotXidMap</code>
<code>stats_reset</code>	<code>timestamp with time zone</code>	Time at which these statistics were last reset
<code>CSNXidMap_last_trim</code>	<code>timestamp with time zone</code>	Shows the last time when the <code>shardman.trim_csnxid_map()</code> function was called.

To reset CSN-related statistics, call the `pg_stat_reset_shared` function with the only text argument equal to `csn`.

Note

Postgres Pro Shardman functionality related to CSN snapshots is work in progress. So anticipate changes to the corresponding views in future releases.

27.2.27.2. shardman.pg_indoubt_xacts

The view `shardman.pg_indoubt_xacts` displays information about transactions that are currently in the `InDoubt` state. An entry is removed when the transaction state changes.

Table 27.38. shardman.pg_indoubt_xacts Columns

Name	Type	Description
xid	xid	Transaction ID of a transaction in the InDoubt state
duration_msec	bigint	Time the transaction was in the InDoubt state, in milliseconds

When the `shardman.pg_indoubt_xacts` view is accessed, the internal transaction manager data structures are momentarily locked, and a copy is made for the view to display. This ensures that the view produces a consistent set of results, while not blocking normal operations longer than necessary. Nonetheless there could be some impact on database performance if this view is frequently accessed.

27.2.27.3. shardman.pg_stat_xact_time

The `shardman.pg_stat_xact_time` view shows statistics for the time spent on a transaction. The columns of the view are shown in [Table 27.39](#).

Table 27.39. shardman.pg_stat_xact_time Columns

Name	Type	Description
overall_committed_xact_time	bigint	Overall time spent for the committed transactions
overall_aborted_xact_time	bigint	Overall time spent for the aborted transactions
overall_commit_time	bigint	Overall time spent for the committing transactions
local_commit_time	bigint	Overall time spent for writing to WAL for all the committed transactions
global_commit_time	bigint	Overall time spent for the distributed queries sending messages about transaction statuses for all the committed transactions
overall_abort_time	bigint	Overall time spent for aborting transactions
local_abort_time	bigint	Overall time spent for writing to WAL for all the aborted transactions
global_abort_time	bigint	Overall time spent for the distributed queries sending messages about transaction statuses for all the aborted transactions
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.27.4. shardman.oldest_csn

The `shardman.oldest_csn` view has one row showing tuple `csn`, `xid`, and `rgid` containing CSN and XID of the oldest transaction in the cluster along with transaction's replication group number.

27.2.27.5. shardman.pg_stat_monitor

The `shardman.pg_stat_monitor` view has one row showing metrics of the Postgres Pro Shardman monitor. The view data is based on the [cumulative statistics system](#). The columns of the view are shown in [Table 27.40](#).

Table 27.40. shardman.pg_stat_monitor Columns

Name	Type	Description
resolved_deadlocks	bigint	Number of resolved distributed deadlocks
aborted_xacts	bigint	Number of aborted outdated prepared transactions
committed_xacts	bigint	Number of committed outdated prepared transactions
errors	bigint	Number of Postgres Pro Shardman monitor errors
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.27.6. shardman.pg_stat_netusage

The `shardman.pg_stat_netusage` view has one row showing the cumulative network traffic between Postgres Pro Shardman cluster nodes. The view data is based on the [cumulative statistics system](#). The columns of the view are shown in [Table 27.41](#).

Table 27.41. shardman.pg_stat_netusage Columns

Name	Type	Description
netusage_recv_bytes	numeric	Total number of bytes received from other nodes through the network by each Postgres Pro Shardman cluster node
netusage_sent_bytes	numeric	Total number of bytes sent to other nodes through the network by each Postgres Pro Shardman cluster node
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.27.7. shardman.pg_stat_foreign_stat_bytes

The `shardman.pg_stat_foreign_stat_bytes` view shows the amount of statistics for foreign relations transferred over the network between Postgres Pro Shardman cluster nodes. The view data is based on the [cumulative statistics system](#). The columns of the view are shown in [Table 27.42](#).

Table 27.42. shardman.pg_stat_foreign_stat_bytes Columns

Name	Type	Description
foreign_stat_recv_bytes	bigint	Total number of bytes of the statistics for the foreign relations received from other nodes through the network by this node
stats_reset	timestamp with time zone	Time at which these statistics were last reset

27.2.28. Statistics Functions

Other ways of looking at the statistics can be set up by writing queries that use the same underlying statistics access functions used by the standard views shown above. For details such as the functions' names, consult the definitions of the standard views. (For example, in `psql` you could issue `\d+ pg_stat_activity`.) The access functions for per-database statistics take a database OID as an argument to

identify which database to report on. The per-table and per-index functions take a table or index OID. The functions for per-function statistics take a function OID. Note that only tables, indexes, and functions in the current database can be seen with these functions.

Additional functions related to the cumulative statistics system are listed in [Table 27.43](#).

Table 27.43. Additional Statistics Functions

Function	Description
<code>pg_backend_pid () → integer</code>	Returns the process ID of the server process attached to the current session.
<code>pg_stat_get_activity (integer) → setof record</code>	Returns a record of information about the backend with the specified process ID, or one record for each active backend in the system if <code>NULL</code> is specified. The fields returned are a subset of those in the <code>pg_stat_activity</code> view.
<code>pg_stat_get_snapshot_timestamp () → timestamp with time zone</code>	Returns the timestamp of the current statistics snapshot, or <code>NULL</code> if no statistics snapshot has been taken. A snapshot is taken the first time cumulative statistics are accessed in a transaction if <code>stats_fetch_consistency</code> is set to <code>snapshot</code> .
<code>pg_stat_get_xact_blocks_fetched (oid) → bigint</code>	Returns the number of block read requests for table or index, in the current transaction. This number minus <code>pg_stat_get_xact_blocks_hit</code> gives the number of kernel <code>read()</code> calls; the number of actual physical reads is usually lower due to kernel-level buffering.
<code>pg_stat_get_xact_blocks_hit (oid) → bigint</code>	Returns the number of block read requests for table or index, in the current transaction, found in cache (not triggering kernel <code>read()</code> calls).
<code>pg_stat_clear_snapshot () → void</code>	Discards the current statistics snapshot or cached information.
<code>pg_stat_reset () → void</code>	Resets all statistics counters for the current database to zero. This function is restricted to superusers by default, but other users can be granted <code>EXECUTE</code> to run the function.
<code>pg_stat_reset_shared ([target text DEFAULT NULL]) → void</code>	<p>Resets some cluster-wide statistics counters to zero, depending on the argument. <i>target</i> can be:</p> <ul style="list-style-type: none"> • <code>archiver</code>: Reset all the counters shown in the <code>pg_stat_archiver</code> view. • <code>bgwriter</code>: Reset all the counters shown in the <code>pg_stat_bgwriter</code> view. • <code>checkpointer</code>: Reset all the counters shown in the <code>pg_stat_checkpointer</code> view. • <code>io</code>: Reset all the counters shown in the <code>pg_stat_io</code> view. • <code>recovery_prefetch</code>: Reset all the counters shown in the <code>pg_stat_recovery_prefetch</code> view. • <code>slru</code>: Reset all the counters shown in the <code>pg_stat_slru</code> view. • <code>wal</code>: Reset all the counters shown in the <code>pg_stat_wal</code> view. • <code>NULL</code> or not specified: All the counters from the views listed above are reset. <p>This function is restricted to superusers by default, but other users can be granted <code>EXECUTE</code> to run the function.</p>
<code>pg_stat_reset_single_table_counters (oid) → void</code>	

Function	Description
	Resets statistics for a single table or index in the current database or shared across all databases in the cluster to zero. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_stat_reset_single_function_counters (oid) → void</code>	Resets statistics for a single function in the current database to zero. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_stat_reset_slru ([target text DEFAULT NULL]) → void</code>	Resets statistics to zero for a single SLRU cache, or for all SLRUs in the cluster. If <i>target</i> is NULL or is not specified, all the counters shown in the <code>pg_stat_slru</code> view for all SLRU caches are reset. The argument can be one of <code>commit_timestamp</code> , <code>multixact_member</code> , <code>multixact_offset</code> , <code>notify</code> , <code>serializable</code> , <code>subtransaction</code> , or <code>transaction</code> to reset the counters for only that entry. If the argument is other (or indeed, any unrecognized name), then the counters for all other SLRU caches, such as extension-defined caches, are reset. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_stat_reset_replication_slot (text) → void</code>	Resets statistics of the replication slot defined by the argument. If the argument is NULL, resets statistics for all the replication slots. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.
<code>pg_stat_reset_subscription_stats (oid) → void</code>	Resets statistics for a single subscription shown in the <code>pg_stat_subscription_stats</code> view to zero. If the argument is NULL, reset statistics for all subscriptions. This function is restricted to superusers by default, but other users can be granted EXECUTE to run the function.

Warning

Using `pg_stat_reset()` also resets counters that autovacuum uses to determine when to trigger a vacuum or an analyze. Resetting these counters can cause autovacuum to not perform necessary work, which can cause problems such as table bloat or out-dated table statistics. A database-wide ANALYZE is recommended after the statistics have been reset.

`pg_stat_get_activity`, the underlying function of the `pg_stat_activity` view, returns a set of records containing all the available information about each backend process. Sometimes it may be more convenient to obtain just a subset of this information. In such cases, another set of per-backend statistics access functions can be used; these are shown in [Table 27.44](#). These access functions use the session's backend ID number, which is a small integer (≥ 0) that is distinct from the backend ID of any concurrent session, although a session's ID can be recycled as soon as it exits. The backend ID is used, among other things, to identify the session's temporary schema if it has one. The function `pg_stat_get_backend_idset` provides a convenient way to list all the active backends' ID numbers for invoking these functions. For example, to show the PIDs and current queries of all backends:

```
SELECT pg_stat_get_backend_pid(backendid) AS pid,
       pg_stat_get_backend_activity(backendid) AS query
FROM pg_stat_get_backend_idset() AS backendid;
```

Table 27.44. Per-Backend Statistics Functions

Function	Description
<code>pg_stat_get_backend_activity (integer) → text</code>	Returns the text of this backend's most recent query.
<code>pg_stat_get_backend_activity_start (integer) → timestamp with time zone</code>	Returns the time when the backend's most recent query was started.
<code>pg_stat_get_backend_client_addr (integer) → inet</code>	Returns the IP address of the client connected to this backend.
<code>pg_stat_get_backend_client_port (integer) → integer</code>	Returns the TCP port number that the client is using for communication.
<code>pg_stat_get_backend_dbid (integer) → oid</code>	Returns the OID of the database this backend is connected to.
<code>pg_stat_get_backend_idset () → setof integer</code>	Returns the set of currently active backend ID numbers.
<code>pg_stat_get_backend_pid (integer) → integer</code>	Returns the process ID of this backend.
<code>pg_stat_get_backend_start (integer) → timestamp with time zone</code>	Returns the time when this process was started.
<code>pg_stat_get_backend_subxact (integer) → record</code>	Returns a record of information about the subtransactions of the backend with the specified ID. The fields returned are <i>subxact_count</i> , which is the number of subtransactions in the backend's subtransaction cache, and <i>subxact_overflow</i> , which indicates whether the backend's subtransaction cache is overflowed or not.
<code>pg_stat_get_backend_userid (integer) → oid</code>	Returns the OID of the user logged into this backend.
<code>pg_stat_get_backend_wait_event (integer) → text</code>	Returns the wait event name if this backend is currently waiting, otherwise NULL. See Table 27.5 through Table 27.13 .
<code>pg_stat_get_backend_wait_event_type (integer) → text</code>	Returns the wait event type name if this backend is currently waiting, otherwise NULL. See Table 27.4 for details.
<code>pg_stat_get_backend_xact_start (integer) → timestamp with time zone</code>	Returns the time when the backend's current transaction was started.
<code>pgpro_stat_get_wal_activity (integer) → bigint</code>	Returns the size of WAL files generated by a backend.

27.3. Viewing Locks

Another useful tool for monitoring database activity is the `pg_locks` system table. It allows the database administrator to view information about the outstanding locks in the lock manager. For example, this capability can be used to:

- View all the locks currently outstanding, all the locks on relations in a particular database, all the locks on a particular relation, or all the locks held by a particular Postgres Pro Shardman session.
- Determine the relation in the current database with the most ungranted locks (which might be a source of contention among database clients).

- Determine the effect of lock contention on overall database performance, as well as the extent to which contention varies with overall database traffic.

Details of the `pg_locks` view appear in [Section 55.12](#). For more information on locking and managing concurrency with Postgres Pro Shardman, refer to [Chapter 13](#).

27.4. Progress Reporting

Postgres Pro Shardman has the ability to report the progress of certain commands during command execution. Currently, the only commands which support progress reporting are `ANALYZE`, `CLUSTER`, `CREATE INDEX`, `VACUUM`, `COPY`, and [BASE_BACKUP](#). This may be expanded in the future.

27.4.1. ANALYZE Progress Reporting

Whenever `ANALYZE` is running, the `pg_stat_progress_analyze` view will contain a row for each backend that is currently running that command. The tables below describe the information that will be reported and provide information about how to interpret it.

Table 27.45. `pg_stat_progress_analyze` View

Column	Type	Description
<code>pid</code>	<code>integer</code>	Process ID of backend.
<code>datid</code>	<code>oid</code>	OID of the database to which this backend is connected.
<code>datname</code>	<code>name</code>	Name of the database to which this backend is connected.
<code>relid</code>	<code>oid</code>	OID of the table being analyzed.
<code>phase</code>	<code>text</code>	Current processing phase. See Table 27.46 .
<code>sample_blks_total</code>	<code>bigint</code>	Total number of heap blocks that will be sampled.
<code>sample_blks_scanned</code>	<code>bigint</code>	Number of heap blocks scanned.
<code>ext_stats_total</code>	<code>bigint</code>	Number of extended statistics.
<code>ext_stats_computed</code>	<code>bigint</code>	Number of extended statistics computed. This counter only advances when the phase is computing extended statistics.
<code>child_tables_total</code>	<code>bigint</code>	Number of child tables.
<code>child_tables_done</code>	<code>bigint</code>	Number of child tables scanned. This counter only advances when the phase is acquiring inherited sample rows.
<code>current_child_table_relid</code>	<code>oid</code>	OID of the child table currently being scanned. This field is only valid when the phase is acquiring inherited sample rows.

Table 27.46. ANALYZE Phases

Phase	Description
<code>initializing</code>	The command is preparing to begin scanning the heap. This phase is expected to be very brief.

Phase	Description
acquiring sample rows	The command is currently scanning the table given by <code>relid</code> to obtain sample rows.
acquiring inherited sample rows	The command is currently scanning child tables to obtain sample rows. Columns <code>child_tables_total</code> , <code>child_tables_done</code> , and <code>current_child_table_relid</code> contain the progress information for this phase.
computing statistics	The command is computing statistics from the sample rows obtained during the table scan.
computing extended statistics	The command is computing extended statistics from the sample rows obtained during the table scan.
finalizing analyze	The command is updating <code>pg_class</code> . When this phase is completed, <code>ANALYZE</code> will end.

Note

Note that when `ANALYZE` is run on a partitioned table, all of its partitions are also recursively analyzed. In that case, `ANALYZE` progress is reported first for the parent table, whereby its inheritance statistics are collected, followed by that for each partition.

27.4.2. CLUSTER Progress Reporting

Whenever `CLUSTER` or `VACUUM FULL` is running, the `pg_stat_progress_cluster` view will contain a row for each backend that is currently running either command. The tables below describe the information that will be reported and provide information about how to interpret it.

Table 27.47. `pg_stat_progress_cluster` View

Column	Type	Description
<code>pid</code>	integer	Process ID of backend.
<code>datid</code>	oid	OID of the database to which this backend is connected.
<code>datname</code>	name	Name of the database to which this backend is connected.
<code>relid</code>	oid	OID of the table being clustered.
<code>command</code>	text	The command that is running. Either <code>CLUSTER</code> or <code>VACUUM FULL</code> .
<code>phase</code>	text	Current processing phase. See Table 27.48 .
<code>cluster_index_relid</code>	oid	If the table is being scanned using an index, this is the OID of the index being used; otherwise, it is zero.
<code>heap_tuples_scanned</code>	bigint	Number of heap tuples scanned. This counter only advances when the phase is <code>seq scanning heap</code> , <code>index scanning heap</code> or writing new heap.
<code>heap_tuples_written</code>	bigint	Number of heap tuples written. This counter only advances when the phase is <code>seq scanning heap</code> , <code>index scanning heap</code> or writing new heap.

Column Type	Description
heap_blks_total bigint	Total number of heap blocks in the table. This number is reported as of the beginning of seq scanning heap.
heap_blks_scanned bigint	Number of heap blocks scanned. This counter only advances when the phase is seq scanning heap.
index_rebuild_count bigint	Number of indexes rebuilt. This counter only advances when the phase is rebuilding index.

Table 27.48. CLUSTER and VACUUM FULL Phases

Phase	Description
initializing	The command is preparing to begin scanning the heap. This phase is expected to be very brief.
seq scanning heap	The command is currently scanning the table using a sequential scan.
index scanning heap	CLUSTER is currently scanning the table using an index scan.
sorting tuples	CLUSTER is currently sorting tuples.
writing new heap	CLUSTER is currently writing the new heap.
swapping relation files	The command is currently swapping newly-built files into place.
rebuilding index	The command is currently rebuilding an index.
performing final cleanup	The command is performing final cleanup. When this phase is completed, CLUSTER or VACUUM FULL will end.

27.4.3. COPY Progress Reporting

Whenever COPY is running, the `pg_stat_progress_copy` view will contain one row for each backend that is currently running a COPY command. The table below describes the information that will be reported and provides information about how to interpret it.

Table 27.49. pg_stat_progress_copy View

Column Type	Description
pid integer	Process ID of backend.
datid oid	OID of the database to which this backend is connected.
datname name	Name of the database to which this backend is connected.
relid oid	OID of the table on which the COPY command is executed. It is set to 0 if copying from a SELECT query.
command text	The command that is running: COPY FROM, or COPY TO.
type text	The I/O type that the data is read from or written to: FILE, PROGRAM, PIPE (for COPY FROM STDIN and COPY TO STDOUT), or CALLBACK (used for example during the initial table synchronization in logical replication).
bytes_processed bigint	

Column Type	Description
	Number of bytes already processed by <code>COPY</code> command.
<code>bytes_total bigint</code>	Size of source file for <code>COPY FROM</code> command in bytes. It is set to 0 if not available.
<code>tuples_processed bigint</code>	Number of tuples already processed by <code>COPY</code> command.
<code>tuples_excluded bigint</code>	Number of tuples not processed because they were excluded by the <code>WHERE</code> clause of the <code>COPY</code> command.
<code>tuples_skipped bigint</code>	Number of tuples skipped because they contain malformed data. This counter only advances when a value other than <code>stop</code> is specified to the <code>ON_ERROR</code> option.

27.4.4. CREATE INDEX Progress Reporting

Whenever `CREATE INDEX` or `REINDEX` is running, the `pg_stat_progress_create_index` view will contain one row for each backend that is currently creating indexes. The tables below describe the information that will be reported and provide information about how to interpret it.

Table 27.50. `pg_stat_progress_create_index` View

Column Type	Description
<code>pid integer</code>	Process ID of the backend creating indexes.
<code>datid oid</code>	OID of the database to which this backend is connected.
<code>datname name</code>	Name of the database to which this backend is connected.
<code>relid oid</code>	OID of the table on which the index is being created.
<code>index_relid oid</code>	OID of the index being created or reindexed. During a non-concurrent <code>CREATE INDEX</code> , this is 0.
<code>command text</code>	Specific command type: <code>CREATE INDEX</code> , <code>CREATE INDEX CONCURRENTLY</code> , <code>REINDEX</code> , or <code>REINDEX CONCURRENTLY</code> .
<code>phase text</code>	Current processing phase of index creation. See Table 27.51 .
<code>lockers_total bigint</code>	Total number of lockers to wait for, when applicable.
<code>lockers_done bigint</code>	Number of lockers already waited for.
<code>current_locker_pid bigint</code>	Process ID of the locker currently being waited for.
<code>blocks_total bigint</code>	Total number of blocks to be processed in the current phase.
<code>blocks_done bigint</code>	Number of blocks already processed in the current phase.
<code>tuples_total bigint</code>	

Column Type	Description
	Total number of tuples to be processed in the current phase.
tuples_done bigint	Number of tuples already processed in the current phase.
partitions_total bigint	Total number of partitions on which the index is to be created or attached, including both direct and indirect partitions. 0 during a REINDEX, or when the index is not partitioned.
partitions_done bigint	Number of partitions on which the index has already been created or attached, including both direct and indirect partitions. 0 during a REINDEX, or when the index is not partitioned.

Table 27.51. CREATE INDEX Phases

Phase	Description
initializing	CREATE INDEX or REINDEX is preparing to create the index. This phase is expected to be very brief.
waiting for writers before build	CREATE INDEX CONCURRENTLY or REINDEX CONCURRENTLY is waiting for transactions with write locks that can potentially see the table to finish. This phase is skipped when not in concurrent mode. Columns lockers_total , lockers_done and current_locker_pid contain the progress information for this phase.
building index	The index is being built by the access method-specific code. In this phase, access methods that support progress reporting fill in their own progress data, and the subphase is indicated in this column. Typically, blocks_total and blocks_done will contain progress data, as well as potentially tuples_total and tuples_done .
waiting for writers before validation	CREATE INDEX CONCURRENTLY or REINDEX CONCURRENTLY is waiting for transactions with write locks that can potentially write into the table to finish. This phase is skipped when not in concurrent mode. Columns lockers_total , lockers_done and current_locker_pid contain the progress information for this phase.
index validation: scanning index	CREATE INDEX CONCURRENTLY is scanning the index searching for tuples that need to be validated. This phase is skipped when not in concurrent mode. Columns blocks_total (set to the total size of the index) and blocks_done contain the progress information for this phase.
index validation: sorting tuples	CREATE INDEX CONCURRENTLY is sorting the output of the index scanning phase.
index validation: scanning table	CREATE INDEX CONCURRENTLY is scanning the table to validate the index tuples collected in the previous two phases. This phase is skipped when not in concurrent mode. Columns blocks_total (set to the total size of the table) and blocks_done contain the progress information for this phase.
waiting for old snapshots	CREATE INDEX CONCURRENTLY or REINDEX CONCURRENTLY is waiting for transactions that can potentially see the table to release their snapshots. This phase is skipped when not in concurrent mode. Columns lockers_total , lockers_done and current_locker_pid contain the progress information for this phase.
waiting for readers before marking dead	REINDEX CONCURRENTLY is waiting for transactions with read locks on the table to finish, before marking the old index dead. This phase is skipped when not in concurrent mode. Columns lock-

Phase	Description
	ers_total , lockers_done and current_locker_pid contain the progress information for this phase.
waiting for readers before dropping	REINDEX CONCURRENTLY is waiting for transactions with read locks on the table to finish, before dropping the old index. This phase is skipped when not in concurrent mode. Columns lockers_total , lockers_done and current_locker_pid contain the progress information for this phase.

27.4.5. VACUUM Progress Reporting

Whenever `VACUUM` is running, the `pg_stat_progress_vacuum` view will contain one row for each backend (including autovacuum worker processes) that is currently vacuuming. The tables below describe the information that will be reported and provide information about how to interpret it. Progress for `VACUUM FULL` commands is reported via `pg_stat_progress_cluster` because both `VACUUM FULL` and `CLUSTER` rewrite the table, while regular `VACUUM` only modifies it in place. See [Section 27.4.2](#).

Table 27.52. `pg_stat_progress_vacuum` View

Column Type	Description
pid integer	Process ID of backend.
datid oid	OID of the database to which this backend is connected.
datname name	Name of the database to which this backend is connected.
relid oid	OID of the table being vacuumed.
phase text	Current processing phase of vacuum. See Table 27.53 .
heap_blks_total bigint	Total number of heap blocks in the table. This number is reported as of the beginning of the scan; blocks added later will not be (and need not be) visited by this <code>VACUUM</code> .
heap_blks_scanned bigint	Number of heap blocks scanned. Because the visibility map is used to optimize scans, some blocks will be skipped without inspection; skipped blocks are included in this total, so that this number will eventually become equal to <code>heap_blks_total</code> when the vacuum is complete. This counter only advances when the phase is scanning heap.
heap_blks_vacuumed bigint	Number of heap blocks vacuumed. Unless the table has no indexes, this counter only advances when the phase is vacuuming heap. Blocks that contain no dead tuples are skipped, so the counter may sometimes skip forward in large increments.
index_vacuum_count bigint	Number of completed index vacuum cycles.
max_dead_tuple_bytes bigint	Amount of dead tuple data that we can store before needing to perform an index vacuum cycle, based on maintenance_work_mem .
dead_tuple_bytes bigint	Amount of dead tuple data collected since the last index vacuum cycle.
num_dead_item_ids bigint	Number of dead item identifiers collected since the last index vacuum cycle.

Column Type	Description
<code>indexes_total bigint</code>	Total number of indexes that will be vacuumed or cleaned up. This number is reported at the beginning of the <code>vacuuming indexes</code> phase or the <code>cleaning up indexes</code> phase.
<code>indexes_processed bigint</code>	Number of indexes processed. This counter only advances when the phase is <code>vacuuming indexes</code> or <code>cleaning up indexes</code> .

Table 27.53. VACUUM Phases

Phase	Description
<code>initializing</code>	VACUUM is preparing to begin scanning the heap. This phase is expected to be very brief.
<code>scanning heap</code>	VACUUM is currently scanning the heap. It will prune and defragment each page if required, and possibly perform freezing activity. The <code>heap_blks_scanned</code> column can be used to monitor the progress of the scan.
<code>vacuuming indexes</code>	VACUUM is currently vacuuming the indexes. If a table has any indexes, this will happen at least once per vacuum, after the heap has been completely scanned. It may happen multiple times per vacuum if <code>maintenance_work_mem</code> (or, in the case of autovacuum, <code>autovacuum_work_mem</code> if set) is insufficient to store the number of dead tuples found.
<code>vacuuming heap</code>	VACUUM is currently vacuuming the heap. Vacuuming the heap is distinct from scanning the heap, and occurs after each instance of vacuuming indexes. If <code>heap_blks_scanned</code> is less than <code>heap_blks_total</code> , the system will return to scanning the heap after this phase is completed; otherwise, it will begin cleaning up indexes after this phase is completed.
<code>cleaning up indexes</code>	VACUUM is currently cleaning up indexes. This occurs after the heap has been completely scanned and all vacuuming of the indexes and the heap has been completed.
<code>truncating heap</code>	VACUUM is currently truncating the heap so as to return empty pages at the end of the relation to the operating system. This occurs after cleaning up indexes.
<code>performing final cleanup</code>	VACUUM is performing final cleanup. During this phase, VACUUM will vacuum the free space map, update statistics in <code>pg_class</code> , and report statistics to the cumulative statistics system. When this phase is completed, VACUUM will end.

27.4.6. Base Backup Progress Reporting

Whenever an application like `pg_basebackup` is taking a base backup, the `pg_stat_progress_basebackup` view will contain a row for each WAL sender process that is currently running the `BASE_BACKUP` replication command and streaming the backup. The tables below describe the information that will be reported and provide information about how to interpret it.

Table 27.54. `pg_stat_progress_basebackup` View

Column Type	Description
<code>pid integer</code>	Process ID of a WAL sender process.

Column Type	Description
phase text	Current processing phase. See Table 27.55 .
backup_total bigint	Total amount of data that will be streamed. This is estimated and reported as of the beginning of streaming database files phase. Note that this is only an approximation since the database may change during streaming database files phase and WAL log may be included in the backup later. This is always the same value as backup_streamed once the amount of data streamed exceeds the estimated total size. If the estimation is disabled in pg_basebackup (i.e., --no-estimate-size option is specified), this is NULL.
backup_streamed bigint	Amount of data streamed. This counter only advances when the phase is streaming database files or transferring wal files.
tablespaces_total bigint	Total number of tablespaces that will be streamed.
tablespaces_streamed bigint	Number of tablespaces streamed. This counter only advances when the phase is streaming database files.

Table 27.55. Base Backup Phases

Phase	Description
initializing	The WAL sender process is preparing to begin the backup. This phase is expected to be very brief.
waiting for checkpoint to finish	The WAL sender process is currently performing pg_backup_start to prepare to take a base backup, and waiting for the start-of-backup checkpoint to finish.
estimating backup size	The WAL sender process is currently estimating the total amount of database files that will be streamed as a base backup.
streaming database files	The WAL sender process is currently streaming database files as a base backup.
waiting for wal archiving to finish	The WAL sender process is currently performing pg_backup_stop to finish the backup, and waiting for all the WAL files required for the base backup to be successfully archived. If either --wal-method=none or --wal-method=stream is specified in pg_basebackup, the backup will end when this phase is completed.
transferring wal files	The WAL sender process is currently transferring all WAL logs generated during the backup. This phase occurs after waiting for wal archiving to finish phase if --wal-method=fetch is specified in pg_basebackup. The backup will end when this phase is completed.

27.5. Monitoring Disk Usage

This section discusses how to monitor the disk usage of a Postgres Pro Shardman database system.

27.5.1. Determining Disk Usage

Each table has a primary heap disk file where most of the data is stored. If the table has any columns with potentially-wide values, there also might be a TOAST file associated with the table, which is used to store values too wide to fit comfortably in the main table (see [Section 66.2](#)). There will be one valid index on the TOAST table, if present. There also might be indexes associated with the base table. Each

table and index is stored in a separate disk file — possibly more than one file, if the file would exceed one gigabyte. Naming conventions for these files are described in [Section 66.1](#).

You can monitor disk space in three ways: using the SQL functions listed in [Table 9.100](#), using the [oid2name](#) module, or using manual inspection of the system catalogs. The SQL functions are the easiest to use and are generally recommended. The remainder of this section shows how to do it by inspection of the system catalogs.

Using `psql` on a recently vacuumed or analyzed database, you can issue queries to see the disk usage of any table:

```
SELECT pg_relation_filepath(oid), relpages FROM pg_class WHERE relname = 'customer';
```

```
pg_relation_filepath | relpages
-----+-----
base/16384/16806    |        60
(1 row)
```

Each page is typically 8 kilobytes. (Remember, `relpages` is only updated by `VACUUM`, `ANALYZE`, and a few DDL commands such as `CREATE INDEX`.) The file path name is of interest if you want to examine the table's disk file directly.

To show the space used by TOAST tables, use a query like the following:

```
SELECT relname, relpages
FROM pg_class,
     (SELECT reltoastrelid
      FROM pg_class
      WHERE relname = 'customer') AS ss
WHERE oid = ss.reltoastrelid OR
      oid = (SELECT indexrelid
            FROM pg_index
            WHERE indrelid = ss.reltoastrelid)
ORDER BY relname;
```

```
relname          | relpages
-----+-----
pg_toast_16806    |         0
pg_toast_16806_index |         1
```

You can easily display index sizes, too:

```
SELECT c2.relname, c2.relpages
FROM pg_class c, pg_class c2, pg_index i
WHERE c.relname = 'customer' AND
      c.oid = i.indrelid AND
      c2.oid = i.indexrelid
ORDER BY c2.relname;
```

```
relname          | relpages
-----+-----
customer_id_index |        26
```

It is easy to find your largest tables and indexes using this information:

```
SELECT relname, relpages
FROM pg_class
ORDER BY relpages DESC;
```

```
relname          | relpages
-----+-----
```

bigtable		3290
customer		3144

27.5.2. Disk Full Failure

The most important disk monitoring task of a database administrator is to make sure the disk doesn't become full. A filled data disk will not result in data corruption, but it might prevent useful activity from occurring. If the disk holding the WAL files grows full, database server panic and consequent shutdown might occur.

If you cannot free up additional space on the disk by deleting other things, you can move some of the database files to other file systems by making use of tablespaces. See [Section 22.7](#) for more information about that.

Tip

Some file systems perform badly when they are almost full, so do not wait until the disk is completely full to take action.

If your system supports per-user disk quotas, then the database will naturally be subject to whatever quota is placed on the user the server runs as. Exceeding the quota will have the same bad effects as running out of disk space entirely.

Chapter 28. Reliability and the Write-Ahead Log

This chapter explains how to control the reliability of Postgres Pro, including details about the Write-Ahead Log.

28.1. Reliability

Reliability is an important property of any serious database system, and Postgres Pro does everything possible to guarantee reliable operation. One aspect of reliable operation is that all data recorded by a committed transaction should be stored in a nonvolatile area that is safe from power loss, operating system failure, and hardware failure (except failure of the nonvolatile area itself, of course). Successfully writing the data to the computer's permanent storage (disk drive or equivalent) ordinarily meets this requirement. In fact, even if a computer is fatally damaged, if the disk drives survive they can be moved to another computer with similar hardware and all committed transactions will remain intact.

While forcing data to the disk platters periodically might seem like a simple operation, it is not. Because disk drives are dramatically slower than main memory and CPUs, several layers of caching exist between the computer's main memory and the disk platters. First, there is the operating system's buffer cache, which caches frequently requested disk blocks and combines disk writes. Fortunately, all operating systems give applications a way to force writes from the buffer cache to disk, and Postgres Pro uses those features. (See the [wal_sync_method](#) parameter to adjust how this is done.)

Next, there might be a cache in the disk drive controller; this is particularly common on RAID controller cards. Some of these caches are *write-through*, meaning writes are sent to the drive as soon as they arrive. Others are *write-back*, meaning data is sent to the drive at some later time. Such caches can be a reliability hazard because the memory in the disk controller cache is volatile, and will lose its contents in a power failure. Better controller cards have *battery-backup units* (BBUs), meaning the card has a battery that maintains power to the cache in case of system power loss. After power is restored the data will be written to the disk drives.

And finally, most disk drives have caches. Some are write-through while some are write-back, and the same concerns about data loss exist for write-back drive caches as for disk controller caches. Consumer-grade IDE and SATA drives are particularly likely to have write-back caches that will not survive a power failure. Many solid-state drives (SSD) also have volatile write-back caches.

These caches can typically be disabled; however, the method for doing this varies by operating system and drive type:

- On Linux, IDE and SATA drives can be queried using `hdparm -I`; write caching is enabled if there is a `* next to Write cache`. `hdparm -W 0` can be used to turn off write caching. SCSI drives can be queried using [sdparm](#). Use `sdparm --get=WCE` to check whether the write cache is enabled and `sdparm --clear=WCE` to disable it.
- On FreeBSD, IDE drives can be queried using `camcontrol identify` and write caching turned off using `hw.ata.wc=0` in `/boot/loader.conf`; SCSI drives can be queried using `camcontrol identify`, and the write cache both queried and changed using `sdparm` when available.
- On Solaris, the disk write cache is controlled by `format -e`. (The Solaris ZFS file system is safe with disk write-cache enabled because it issues its own disk cache flush commands.)
- On Windows, if `wal_sync_method` is `open_datasync` (the default), write caching can be disabled by unchecking `My Computer\Open\disk drive\Properties\Hardware\Properties\Policies\Enable write caching on the disk`. Alternatively, set `wal_sync_method` to `fdasync` (NTFS only) or `fsync`, which prevent write caching.
- On macOS, write caching can be prevented by setting `wal_sync_method` to `fsync_writethrough`.

Recent SATA drives (those following ATA-6 or later) offer a drive cache flush command (`FLUSH CACHE EXT`), while SCSI drives have long supported a similar command `SYNCHRONIZE CACHE`. These commands

are not directly accessible to Postgres Pro, but some file systems (e.g., ZFS, ext4) can use them to flush data to the platters on write-back-enabled drives. Unfortunately, such file systems behave suboptimally when combined with battery-backup unit (BBU) disk controllers. In such setups, the `synchronize` command forces all data from the controller cache to the disks, eliminating much of the benefit of the BBU. You can run the [pg_test_fsync](#) program to see if you are affected. If you are affected, the performance benefits of the BBU can be regained by turning off write barriers in the file system or reconfiguring the disk controller, if that is an option. If write barriers are turned off, make sure the battery remains functional; a faulty battery can potentially lead to data loss. Hopefully file system and disk controller designers will eventually address this suboptimal behavior.

When the operating system sends a write request to the storage hardware, there is little it can do to make sure the data has arrived at a truly non-volatile storage area. Rather, it is the administrator's responsibility to make certain that all storage components ensure integrity for both data and file-system metadata. Avoid disk controllers that have non-battery-backed write caches. At the drive level, disable write-back caching if the drive cannot guarantee the data will be written before shutdown. If you use SSDs, be aware that many of these do not honor cache flush commands by default. You can test for reliable I/O subsystem behavior using [diskchecker.pl](#).

Another risk of data loss is posed by the disk platter write operations themselves. Disk platters are divided into sectors, commonly 512 bytes each. Every physical read or write operation processes a whole sector. When a write request arrives at the drive, it might be for some multiple of 512 bytes (Postgres Pro typically writes 8192 bytes, or 16 sectors, at a time), and the process of writing could fail due to power loss at any time, meaning some of the 512-byte sectors were written while others were not. To guard against such failures, Postgres Pro periodically writes full page images to permanent WAL storage *before* modifying the actual page on disk. By doing this, during crash recovery Postgres Pro can restore partially-written pages from WAL. If you have file-system software that prevents partial page writes (e.g., ZFS), you can turn off this page imaging by turning off the [full_page_writes](#) parameter. Battery-Backed Unit (BBU) disk controllers do not prevent partial page writes unless they guarantee that data is written to the BBU as full (8kB) pages.

Postgres Pro also protects against some kinds of data corruption on storage devices that may occur because of hardware errors or media failure over time, such as reading/writing garbage data.

- Each individual record in a WAL file is protected by a CRC-32C (32-bit) check that allows us to tell if record contents are correct. The CRC value is set when we write each WAL record and checked during crash recovery, archive recovery and replication.
- Data pages are not currently checksummed by default, though full page images recorded in WAL records will be protected; see [initdb](#) for details about enabling data checksums.
- Internal data structures such as `pg_xact`, `pg_subtrans`, `pg_multixact`, `pg_serial`, `pg_notify`, `pg_stat`, `pg_snapshots` are not directly checksummed, nor are pages protected by full page writes. However, where such data structures are persistent, WAL records are written that allow recent changes to be accurately rebuilt at crash recovery and those WAL records are protected as discussed above.
- Individual state files in `pg_twophase` are protected by CRC-32C.
- Temporary data files used in larger SQL queries for sorts, materializations and intermediate results are not currently checksummed, nor will WAL records be written for changes to those files.

Postgres Pro does not protect against correctable memory errors and it is assumed you will operate using RAM that uses industry standard Error Correcting Codes (ECC) or better protection.

28.2. Data Checksums

By default, data pages are not protected by checksums, but this can optionally be enabled for a cluster. When enabled, each data page includes a checksum that is updated when the page is written and verified each time the page is read. Only data pages are protected by checksums; internal data structures and temporary files are not.

Checksums can be enabled when the cluster is initialized using [initdb](#). They can also be enabled or disabled at a later time as an offline operation. Data checksums are enabled or disabled at the full cluster level, and cannot be specified individually for databases or tables.

The current state of checksums in the cluster can be verified by viewing the value of the read-only configuration variable [data_checksums](#) by issuing the command `SHOW data_checksums`.

When attempting to recover from page corruptions, it may be necessary to bypass the checksum protection. To do this, temporarily set the configuration parameter [ignore_checksum_failure](#).

28.2.1. Off-line Enabling of Checksums

The [pg_checksums](#) application can be used to enable or disable data checksums, as well as verify checksums, on an offline cluster.

28.3. Write-Ahead Logging (WAL)

Write-Ahead Logging (WAL) is a standard method for ensuring data integrity. A detailed description can be found in most (if not all) books about transaction processing. Briefly, WAL's central concept is that changes to data files (where tables and indexes reside) must be written only after those changes have been logged, that is, after WAL records describing the changes have been flushed to permanent storage. If we follow this procedure, we do not need to flush data pages to disk on every transaction commit, because we know that in the event of a crash we will be able to recover the database using the log: any changes that have not been applied to the data pages can be redone from the WAL records. (This is roll-forward recovery, also known as REDO.)

Tip

Because WAL restores database file contents after a crash, journaled file systems are not necessary for reliable storage of the data files or WAL files. In fact, journaling overhead can reduce performance, especially if journaling causes file system *data* to be flushed to disk. Fortunately, data flushing during journaling can often be disabled with a file system mount option, e.g., `data=writeback` on a Linux ext3 file system. Journaled file systems do improve boot speed after a crash.

Using WAL results in a significantly reduced number of disk writes, because only the WAL file needs to be flushed to disk to guarantee that a transaction is committed, rather than every data file changed by the transaction. The WAL file is written sequentially, and so the cost of syncing the WAL is much less than the cost of flushing the data pages. This is especially true for servers handling many small transactions touching different parts of the data store. Furthermore, when the server is processing many small concurrent transactions, one `fsync` of the WAL file may suffice to commit many transactions.

WAL also makes it possible to support on-line backup and point-in-time recovery, as described in [Section 25.1](#). By archiving the WAL data we can support reverting to any time instant covered by the available WAL data: we simply install a prior physical backup of the database, and replay the WAL just as far as the desired time. What's more, the physical backup doesn't have to be an instantaneous snapshot of the database state — if it is made over some period of time, then replaying the WAL for that period will fix any internal inconsistencies.

28.4. Asynchronous Commit

Asynchronous commit is an option that allows transactions to complete more quickly, at the cost that the most recent transactions may be lost if the database should crash. In many applications this is an acceptable trade-off.

As described in the previous section, transaction commit is normally *synchronous*: the server waits for the transaction's WAL records to be flushed to permanent storage before returning a success indication

to the client. The client is therefore guaranteed that a transaction reported to be committed will be preserved, even in the event of a server crash immediately after. However, for short transactions this delay is a major component of the total transaction time. Selecting asynchronous commit mode means that the server returns success as soon as the transaction is logically completed, before the WAL records it generated have actually made their way to disk. This can provide a significant boost in throughput for small transactions.

Asynchronous commit introduces the risk of data loss. There is a short time window between the report of transaction completion to the client and the time that the transaction is truly committed (that is, it is guaranteed not to be lost if the server crashes). Thus asynchronous commit should not be used if the client will take external actions relying on the assumption that the transaction will be remembered. As an example, a bank would certainly not use asynchronous commit for a transaction recording an ATM's dispensing of cash. But in many scenarios, such as event logging, there is no need for a strong guarantee of this kind.

The risk that is taken by using asynchronous commit is of data loss, not data corruption. If the database should crash, it will recover by replaying WAL up to the last record that was flushed. The database will therefore be restored to a self-consistent state, but any transactions that were not yet flushed to disk will not be reflected in that state. The net effect is therefore loss of the last few transactions. Because the transactions are replayed in commit order, no inconsistency can be introduced — for example, if transaction B made changes relying on the effects of a previous transaction A, it is not possible for A's effects to be lost while B's effects are preserved.

The user can select the commit mode of each transaction, so that it is possible to have both synchronous and asynchronous commit transactions running concurrently. This allows flexible trade-offs between performance and certainty of transaction durability. The commit mode is controlled by the user-settable parameter `synchronous_commit`, which can be changed in any of the ways that a configuration parameter can be set. The mode used for any one transaction depends on the value of `synchronous_commit` when transaction commit begins.

Certain utility commands, for instance `DROP TABLE`, are forced to commit synchronously regardless of the setting of `synchronous_commit`. This is to ensure consistency between the server's file system and the logical state of the database. The commands supporting two-phase commit, such as `PREPARE TRANSACTION`, are also always synchronous.

If the database crashes during the risk window between an asynchronous commit and the writing of the transaction's WAL records, then changes made during that transaction *will* be lost. The duration of the risk window is limited because a background process (the “WAL writer”) flushes unwritten WAL records to disk every `wal_writer_delay` milliseconds. The actual maximum duration of the risk window is three times `wal_writer_delay` because the WAL writer is designed to favor writing whole pages at a time during busy periods.

Caution

An immediate-mode shutdown is equivalent to a server crash, and will therefore cause loss of any unflushed asynchronous commits.

Asynchronous commit provides behavior different from setting `fsync = off`. `fsync` is a server-wide setting that will alter the behavior of all transactions. It disables all logic within Postgres Pro that attempts to synchronize writes to different portions of the database, and therefore a system crash (that is, a hardware or operating system crash, not a failure of Postgres Pro itself) could result in arbitrarily bad corruption of the database state. In many scenarios, asynchronous commit provides most of the performance improvement that could be obtained by turning off `fsync`, but without the risk of data corruption.

`commit_delay` also sounds very similar to asynchronous commit, but it is actually a synchronous commit method (in fact, `commit_delay` is ignored during an asynchronous commit). `commit_delay` causes a delay just before a transaction flushes WAL to disk, in the hope that a single flush executed by one such

transaction can also serve other transactions committing at about the same time. The setting can be thought of as a way of increasing the time window in which transactions can join a group about to participate in a single flush, to amortize the cost of the flush among multiple transactions.

28.5. WAL Configuration

There are several WAL-related configuration parameters that affect database performance. This section explains their use. Consult [Chapter 19](#) for general information about setting server configuration parameters.

Checkpoints are points in the sequence of transactions at which it is guaranteed that the heap and index data files have been updated with all information written before that checkpoint. At checkpoint time, all dirty data pages are flushed to disk and a special checkpoint record is written to the WAL file. (The change records were previously flushed to the WAL files.) In the event of a crash, the crash recovery procedure looks at the latest checkpoint record to determine the point in the WAL (known as the redo record) from which it should start the REDO operation. Any changes made to data files before that point are guaranteed to be already on disk. Hence, after a checkpoint, WAL segments preceding the one containing the redo record are no longer needed and can be recycled or removed. (When WAL archiving is being done, the WAL segments must be archived before being recycled or removed.)

The checkpoint requirement of flushing all dirty data pages to disk can cause a significant I/O load. For this reason, checkpoint activity is throttled so that I/O begins at checkpoint start and completes before the next checkpoint is due to start; this minimizes performance degradation during checkpoints.

The server's checkpoint process automatically performs a checkpoint every so often. A checkpoint is begun every `checkpoint_timeout` seconds, or if `max_wal_size` is about to be exceeded, whichever comes first. The default settings are 5 minutes and 1 GB, respectively. If no WAL has been written since the previous checkpoint, new checkpoints will be skipped even if `checkpoint_timeout` has passed. (If WAL archiving is being used and you want to put a lower limit on how often files are archived in order to bound potential data loss, you should adjust the `archive_timeout` parameter rather than the checkpoint parameters.) It is also possible to force a checkpoint by using the SQL command `CHECKPOINT`.

Reducing `checkpoint_timeout` and/or `max_wal_size` causes checkpoints to occur more often. This allows faster after-crash recovery, since less work will need to be redone. However, one must balance this against the increased cost of flushing dirty data pages more often. If `full_page_writes` is set (as is the default), there is another factor to consider. To ensure data page consistency, the first modification of a data page after each checkpoint results in logging the entire page content. In that case, a smaller checkpoint interval increases the volume of output to the WAL, partially negating the goal of using a smaller interval, and in any case causing more disk I/O.

Checkpoints are fairly expensive, first because they require writing out all currently dirty buffers, and second because they result in extra subsequent WAL traffic as discussed above. It is therefore wise to set the checkpointing parameters high enough so that checkpoints don't happen too often. As a simple sanity check on your checkpointing parameters, you can set the `checkpoint_warning` parameter. If checkpoints happen closer together than `checkpoint_warning` seconds, a message will be output to the server log recommending increasing `max_wal_size`. Occasional appearance of such a message is not cause for alarm, but if it appears often then the checkpoint control parameters should be increased. Bulk operations such as large `COPY` transfers might cause a number of such warnings to appear if you have not set `max_wal_size` high enough.

To avoid flooding the I/O system with a burst of page writes, writing dirty buffers during a checkpoint is spread over a period of time. That period is controlled by `checkpoint_completion_target`, which is given as a fraction of the checkpoint interval (configured by using `checkpoint_timeout`). The I/O rate is adjusted so that the checkpoint finishes when the given fraction of `checkpoint_timeout` seconds have elapsed, or before `max_wal_size` is exceeded, whichever is sooner. With the default value of 0.9, Postgres Pro can be expected to complete each checkpoint a bit before the next scheduled checkpoint (at around 90% of the last checkpoint's duration). This spreads out the I/O as much as possible so that the checkpoint I/O load is consistent throughout the checkpoint interval. The disadvantage of this is that prolonging checkpoints affects recovery time, because more WAL segments will need to be kept

around for possible use in recovery. A user concerned about the amount of time required to recover might wish to reduce `checkpoint_timeout` so that checkpoints occur more frequently but still spread the I/O across the checkpoint interval. Alternatively, `checkpoint_completion_target` could be reduced, but this would result in times of more intense I/O (during the checkpoint) and times of less I/O (after the checkpoint completed but before the next scheduled checkpoint) and therefore is not recommended. Although `checkpoint_completion_target` could be set as high as 1.0, it is typically recommended to set it to no higher than 0.9 (the default) since checkpoints include some other activities besides writing dirty buffers. A setting of 1.0 is quite likely to result in checkpoints not being completed on time, which would result in performance loss due to unexpected variation in the number of WAL segments needed.

On Linux and POSIX platforms `checkpoint_flush_after` allows you to force OS pages written by the checkpoint to be flushed to disk after a configurable number of bytes. Otherwise, these pages may be kept in the OS's page cache, inducing a stall when `fsync` is issued at the end of a checkpoint. This setting will often help to reduce transaction latency, but it also can have an adverse effect on performance; particularly for workloads that are bigger than `shared_buffers`, but smaller than the OS's page cache.

The number of WAL segment files in `pg_wal` directory depends on `min_wal_size`, `max_wal_size` and the amount of WAL generated in previous checkpoint cycles. When old WAL segment files are no longer needed, they are removed or recycled (that is, renamed to become future segments in the numbered sequence). If, due to a short-term peak of WAL output rate, `max_wal_size` is exceeded, the unneeded segment files will be removed until the system gets back under this limit. Below that limit, the system recycles enough WAL files to cover the estimated need until the next checkpoint, and removes the rest. The estimate is based on a moving average of the number of WAL files used in previous checkpoint cycles. The moving average is increased immediately if the actual usage exceeds the estimate, so it accommodates peak usage rather than average usage to some extent. `min_wal_size` puts a minimum on the amount of WAL files recycled for future usage; that much WAL is always recycled for future use, even if the system is idle and the WAL usage estimate suggests that little WAL is needed.

Independently of `max_wal_size`, the most recent `wal_keep_size` megabytes of WAL files plus one additional WAL file are kept at all times. Also, if WAL archiving is used, old segments cannot be removed or recycled until they are archived. If WAL archiving cannot keep up with the pace that WAL is generated, or if `archive_command` or `archive_library` fails repeatedly, old WAL files will accumulate in `pg_wal` until the situation is resolved. A slow or failed standby server that uses a replication slot will have the same effect (see [Section 26.2.6](#)). Similarly, if [WAL summarization](#) is enabled, old segments are kept until they are summarized.

In archive recovery or standby mode, the server periodically performs *restartpoints*, which are similar to checkpoints in normal operation: the server forces all its state to disk, updates the `pg_control` file to indicate that the already-processed WAL data need not be scanned again, and then recycles any old WAL segment files in the `pg_wal` directory. Restartpoints can't be performed more frequently than checkpoints on the primary because restartpoints can only be performed at checkpoint records. A restartpoint can be demanded by a schedule or by an external request. The `restartpoints_timed` counter in the `pg_stat_checkpoint` view counts the first ones while the `restartpoints_req` the second. A restartpoint is triggered by schedule when a checkpoint record is reached if at least `checkpoint_timeout` seconds have passed since the last performed restartpoint or when the previous attempt to perform the restartpoint has failed. In the last case, the next restartpoint will be scheduled in 15 seconds. A restartpoint is triggered by request due to similar reasons like checkpoint but mostly if WAL size is about to exceed `max_wal_size`. However, because of limitations on when a restartpoint can be performed, `max_wal_size` is often exceeded during recovery, by up to one checkpoint cycle's worth of WAL. (`max_wal_size` is never a hard limit anyway, so you should always leave plenty of headroom to avoid running out of disk space.) The `restartpoints_done` counter in the `pg_stat_checkpoint` view counts the restartpoints that have really been performed.

In some cases, when the WAL size on the primary increases quickly, for instance during massive INSERT, the `restartpoints_req` counter on the standby may demonstrate a peak growth. This occurs because requests to create a new restartpoint due to increased XLOG consumption cannot be performed because the safe checkpoint record since the last restartpoint has not yet been replayed on the standby. This behavior is normal and does not lead to an increase in system resource consumption. Only the `restart-`

`points_done` counter among the restartpoint-related ones indicates that noticeable system resources have been spent.

There are two commonly used internal WAL functions: `XLogInsertRecord` and `XLogFlush`. `XLogInsertRecord` is used to place a new record into the WAL buffers in shared memory. If there is no space for the new record, `XLogInsertRecord` will have to write (move to kernel cache) a few filled WAL buffers. This is undesirable because `XLogInsertRecord` is used on every database low level modification (for example, row insertion) at a time when an exclusive lock is held on affected data pages, so the operation needs to be as fast as possible. What is worse, writing WAL buffers might also force the creation of a new WAL segment, which takes even more time. Normally, WAL buffers should be written and flushed by an `XLogFlush` request, which is made, for the most part, at transaction commit time to ensure that transaction records are flushed to permanent storage. On systems with high WAL output, `XLogFlush` requests might not occur often enough to prevent `XLogInsertRecord` from having to do writes. On such systems one should increase the number of WAL buffers by modifying the `wal_buffers` parameter. When `full_page_writes` is set and the system is very busy, setting `wal_buffers` higher will help smooth response times during the period immediately following each checkpoint.

The `commit_delay` parameter defines for how many microseconds a group commit leader process will sleep after acquiring a lock within `XLogFlush`, while group commit followers queue up behind the leader. This delay allows other server processes to add their commit records to the WAL buffers so that all of them will be flushed by the leader's eventual sync operation. No sleep will occur if `fsync` is not enabled, or if fewer than `commit_siblings` other sessions are currently in active transactions; this avoids sleeping when it's unlikely that any other session will commit soon. Note that on some platforms, the resolution of a sleep request is ten milliseconds, so that any nonzero `commit_delay` setting between 1 and 10000 microseconds would have the same effect. Note also that on some platforms, sleep operations may take slightly longer than requested by the parameter.

Since the purpose of `commit_delay` is to allow the cost of each flush operation to be amortized across concurrently committing transactions (potentially at the expense of transaction latency), it is necessary to quantify that cost before the setting can be chosen intelligently. The higher that cost is, the more effective `commit_delay` is expected to be in increasing transaction throughput, up to a point. The `pg_test_fsync` program can be used to measure the average time in microseconds that a single WAL flush operation takes. A value of half of the average time the program reports it takes to flush after a single 8kB write operation is often the most effective setting for `commit_delay`, so this value is recommended as the starting point to use when optimizing for a particular workload. While tuning `commit_delay` is particularly useful when the WAL is stored on high-latency rotating disks, benefits can be significant even on storage media with very fast sync times, such as solid-state drives or RAID arrays with a battery-backed write cache; but this should definitely be tested against a representative workload. Higher values of `commit_siblings` should be used in such cases, whereas smaller `commit_siblings` values are often helpful on higher latency media. Note that it is quite possible that a setting of `commit_delay` that is too high can increase transaction latency by so much that total transaction throughput suffers.

When `commit_delay` is set to zero (the default), it is still possible for a form of group commit to occur, but each group will consist only of sessions that reach the point where they need to flush their commit records during the window in which the previous flush operation (if any) is occurring. At higher client counts a “gangway effect” tends to occur, so that the effects of group commit become significant even when `commit_delay` is zero, and thus explicitly setting `commit_delay` tends to help less. Setting `commit_delay` can only help when (1) there are some concurrently committing transactions, and (2) throughput is limited to some degree by commit rate; but with high rotational latency this setting can be effective in increasing transaction throughput with as few as two clients (that is, a single committing client with one sibling transaction).

The `wal_sync_method` parameter determines how Postgres Pro will ask the kernel to force WAL updates out to disk. All the options should be the same in terms of reliability, with the exception of `fsync_writethrough`, which can sometimes force a flush of the disk cache even when other options do not do so. However, it's quite platform-specific which one will be the fastest. You can test the speeds of different options using the `pg_test_fsync` program. Note that this parameter is irrelevant if `fsync` has been turned off.

Enabling the [wal_debug](#) configuration parameter (provided that Postgres Pro has been compiled with support for it) will result in each `XLogInsertRecord` and `XLogFlush` WAL call being logged to the server log. This option might be replaced by a more general mechanism in the future.

There are two internal functions to write WAL data to disk: `XLogWrite` and `issue_xlog_fsync`. When [track_wal_io_timing](#) is enabled, the total amounts of time `XLogWrite` writes and `issue_xlog_fsync` syncs WAL data to disk are counted as `wal_write_time` and `wal_sync_time` in [pg_stat_wal](#), respectively. `XLogWrite` is normally called by `XLogInsertRecord` (when there is no space for the new record in WAL buffers), `XLogFlush` and the WAL writer, to write WAL buffers to disk and call `issue_xlog_fsync`. `issue_xlog_fsync` is normally called by `XLogWrite` to sync WAL files to disk. If `wal_sync_method` is either `open_datasync` or `open_sync`, a write operation in `XLogWrite` guarantees to sync written WAL data to disk and `issue_xlog_fsync` does nothing. If `wal_sync_method` is either `fdatasync`, `fsync`, or `fsync_writethrough`, the write operation moves WAL buffers to kernel cache and `issue_xlog_fsync` syncs them to disk. Regardless of the setting of `track_wal_io_timing`, the number of times `XLogWrite` writes and `issue_xlog_fsync` syncs WAL data to disk are also counted as `wal_write` and `wal_sync` in [pg_stat_wal](#), respectively.

The [recovery_prefetch](#) parameter can be used to reduce I/O wait times during recovery by instructing the kernel to initiate reads of disk blocks that will soon be needed but are not currently in Postgres Pro's buffer pool. The [maintenance_io_concurrency](#) and [wal_decode_buffer_size](#) settings limit prefetching concurrency and distance, respectively. By default, it is set to `try`, which enables the feature on systems where `posix_fadvise` is available.

28.6. WAL Restoration

Postgres Pro uses WAL to provide protection against some types of data corruption that may occur because of hardware faults. To validate data integrity, all the WAL records are protected by a CRC check.

However, some uncommon transient faults that can be attributed to hardware may cause checksum errors when moving WAL records. As a result, WAL changes would not apply to a replica, and the physical replication would terminate due to a fatal error accompanied by an error message about an incorrect data checksum.

Postgres Pro offers a solution to this type of issues by restoring corrupted WAL data from in-memory WAL buffers. The amount of WAL data kept in memory is specified by the [wal_buffers](#) parameter. The WAL sender process checks CRC values of WAL records before sending them to a replica. If a corrupted record is detected, the WAL sender process tries to restore it from the buffers. For extra protection, there are two copies of the buffers. If the corrupted record could not be restored, the WAL sender process aborts replication with an error. The level of error depends on the value of the [wal_sender_panic_on_crc_error](#) configuration parameter. So if you already have, or suspect you have, such hardware problems, it is recommended to enable additional WAL restoration from WAL buffers by modifying the [wal_sender_check_crc](#) parameter.

28.7. WAL Internals

WAL is automatically enabled; no action is required from the administrator except ensuring that the disk-space requirements for the WAL files are met, and that any necessary tuning is done (see [Section 28.5](#)).

WAL records are appended to the WAL files as each new record is written. The insert position is described by a Log Sequence Number (LSN) that is a byte offset into the WAL, increasing monotonically with each new record. LSN values are returned as the datatype [pg_lsn](#). Values can be compared to calculate the volume of WAL data that separates them, so they are used to measure the progress of replication and recovery.

WAL files are stored in the directory `pg_wal` under the data directory, as a set of segment files, normally each 16 MB in size (but the size can be changed by altering the `--wal-segsize` `initdb` option). Each segment is divided into pages, normally 8 kB each. The log record content is dependent on the type of event that is being logged. Segment files are given ever-increasing numbers as names, starting at

00000000100000000000000001. The numbers do not wrap, but it will take a very, very long time to exhaust the available stock of numbers.

It is advantageous if the WAL is located on a different disk from the main database files. This can be achieved by moving the `pg_wal` directory to another location (while the server is shut down, of course) and creating a symbolic link from the original location in the main data directory to the new location.

The aim of WAL is to ensure that the log is written before database records are altered, but this can be subverted by disk drives that falsely report a successful write to the kernel, when in fact they have only cached the data and not yet stored it on the disk. A power failure in such a situation might lead to irrecoverable data corruption. Administrators should try to ensure that disks holding Postgres Pro's WAL files do not make such false reports. (See [Section 28.1](#).)

After a checkpoint has been made and the WAL flushed, the checkpoint's position is saved in the file `pg_control`. Therefore, at the start of recovery, the server first reads `pg_control` and then the checkpoint record; then it performs the REDO operation by scanning forward from the WAL location indicated in the checkpoint record. Because the entire content of data pages is saved in the WAL on the first page modification after a checkpoint (assuming `full_page_writes` is not disabled), all pages changed since the checkpoint will be restored to a consistent state.

To deal with the case where `pg_control` is corrupt, we should support the possibility of scanning existing WAL segments in reverse order — newest to oldest — in order to find the latest checkpoint. This has not been implemented yet. `pg_control` is small enough (less than one disk page) that it is not subject to partial-write problems, and as of this writing there have been no reports of database failures due solely to the inability to read `pg_control` itself. So while it is theoretically a weak spot, `pg_control` does not seem to be a problem in practice.

Chapter 29. Logical Replication

Logical replication is a method of replicating data objects and their changes, based upon their replication identity (usually a primary key). We use the term logical in contrast to physical replication, which uses exact block addresses and byte-by-byte replication. Postgres Pro supports both mechanisms concurrently, see [Chapter 26](#). Logical replication allows fine-grained control over both data replication and security.

Logical replication uses a *publish* and *subscribe* model with one or more *subscribers* subscribing to one or more *publications* on a *publisher* node. Subscribers pull data from the publications they subscribe to and may subsequently re-publish data to allow cascading replication or more complex configurations.

Logical replication of a table typically starts with taking a snapshot of the data on the publisher database and copying that to the subscriber. Once that is done, the changes on the publisher are sent to the subscriber as they occur in real-time. The subscriber applies the data in the same order as the publisher so that transactional consistency is guaranteed for publications within a single subscription. This method of data replication is sometimes referred to as transactional replication.

The typical use-cases for logical replication are:

- Sending incremental changes in a single database or a subset of a database to subscribers as they occur.
- Firing triggers for individual changes as they arrive on the subscriber.
- Consolidating multiple databases into a single one (for example for analytical purposes).
- Replicating between different major versions of Postgres Pro.
- Replicating between Postgres Pro instances on different platforms (for example Linux to Windows)
- Giving access to replicated data to different groups of users.
- Sharing a subset of the database between multiple databases.

The subscriber database behaves in the same way as any other Postgres Pro instance and can be used as a publisher for other databases by defining its own publications. When the subscriber is treated as read-only by application, there will be no conflicts from a single subscription. On the other hand, if there are other writes done either by an application or by other subscribers to the same set of tables, conflicts can arise.

Note

Postgres Pro supports logical replication to a Postgres Pro Shardman server from servers of different editions: PostgreSQL or Postgres Pro Standard of the same or a previous version.

29.1. Publication

A *publication* can be defined on any physical replication primary. The node where a publication is defined is referred to as *publisher*. A publication is a set of changes generated from a table or a group of tables, and might also be described as a change set or replication set. Each publication exists in only one database.

Publications are different from schemas and do not affect how the table is accessed. Each table can be added to multiple publications if needed. Publications may currently only contain tables and all tables in schema. Objects must be added explicitly, except when a publication is created for `ALL TABLES`.

Publications can choose to limit the changes they produce to any combination of `INSERT`, `UPDATE`, `DELETE`, and `TRUNCATE`, similar to how triggers are fired by particular event types. By default, all operation types

are replicated. These publication specifications apply only for DML operations; they do not affect the initial data synchronization copy. (Row filters have no effect for `TRUNCATE`. See [Section 29.4](#)).

A published table must have a *replica identity* configured in order to be able to replicate `UPDATE` and `DELETE` operations, so that appropriate rows to update or delete can be identified on the subscriber side. By default, this is the primary key, if there is one. Another unique index (with certain additional requirements) can also be set to be the replica identity. If the table does not have any suitable key, then it can be set to replica identity `FULL`, which means the entire row becomes the key. When replica identity `FULL` is specified, indexes can be used on the subscriber side for searching the rows. Candidate indexes must be `btree` or `hash`, non-partial, and the leftmost index field must be a column (not an expression) that references the published table column. These restrictions on the non-unique index properties adhere to some of the restrictions that are enforced for primary keys. If there are no such suitable indexes, the search on the subscriber side can be very inefficient, therefore replica identity `FULL` should only be used as a fallback if no other solution is possible. If a replica identity other than `FULL` is set on the publisher side, a replica identity comprising the same or fewer columns must also be set on the subscriber side. See [REPLICA IDENTITY](#) for details on how to set the replica identity. If a table without a replica identity is added to a publication that replicates `UPDATE` or `DELETE` operations then subsequent `UPDATE` or `DELETE` operations will cause an error on the publisher. `INSERT` operations can proceed regardless of any replica identity.

Every publication can have multiple subscribers.

A publication is created using the [CREATE PUBLICATION](#) command and may later be altered or dropped using corresponding commands.

The individual tables can be added and removed dynamically using [ALTER PUBLICATION](#). Both the `ADD TABLE` and `DROP TABLE` operations are transactional; so the table will start or stop replicating at the correct snapshot once the transaction has committed.

29.2. Subscription

A *subscription* is the downstream side of logical replication. The node where a subscription is defined is referred to as the *subscriber*. A subscription defines the connection to another database and set of publications (one or more) to which it wants to subscribe.

The subscriber database behaves in the same way as any other Postgres Pro instance and can be used as a publisher for other databases by defining its own publications.

A subscriber node may have multiple subscriptions if desired. It is possible to define multiple subscriptions between a single publisher-subscriber pair, in which case care must be taken to ensure that the subscribed publication objects don't overlap.

Each subscription will receive changes via one replication slot (see [Section 26.2.6](#)). Additional replication slots may be required for the initial data synchronization of pre-existing table data and those will be dropped at the end of data synchronization.

A logical replication subscription can be a standby for synchronous replication (see [Section 26.2.8](#)). The standby name is by default the subscription name. An alternative name can be specified as `application_name` in the connection information of the subscription.

Subscriptions are dumped by `pg_dump` if the current user is a superuser. Otherwise a warning is written and subscriptions are skipped, because non-superusers cannot read all subscription information from the `pg_subscription` catalog.

The subscription is added using [CREATE SUBSCRIPTION](#) and can be stopped/resumed at any time using the [ALTER SUBSCRIPTION](#) command and removed using [DROP SUBSCRIPTION](#).

When a subscription is dropped and recreated, the synchronization information is lost. This means that the data has to be resynchronized afterwards.

The schema definitions are not replicated, and the published tables must exist on the subscriber. Only regular tables may be the target of replication. For example, you can't replicate to a view.

The tables are matched between the publisher and the subscriber using the fully qualified table name. Replication to differently-named tables on the subscriber is not supported.

Columns of a table are also matched by name. The order of columns in the subscriber table does not need to match that of the publisher. The data types of the columns do not need to match, as long as the text representation of the data can be converted to the target type. For example, you can replicate from a column of type `integer` to a column of type `bigint`. The target table can also have additional columns not provided by the published table. Any such columns will be filled with the default value as specified in the definition of the target table. However, logical replication in binary format is more restrictive. See the [binary](#) option of `CREATE SUBSCRIPTION` for details.

29.2.1. Replication Slot Management

As mentioned earlier, each (active) subscription receives changes from a replication slot on the remote (publishing) side.

Additional table synchronization slots are normally transient, created internally to perform initial table synchronization and dropped automatically when they are no longer needed. These table synchronization slots have generated names: `"pg_%u_sync_%u_%llu"` (parameters: Subscription *oid*, Table *relid*, system identifier *sysid*)

Normally, the remote replication slot is created automatically when the subscription is created using `CREATE SUBSCRIPTION` and it is dropped automatically when the subscription is dropped using `DROP SUBSCRIPTION`. In some situations, however, it can be useful or necessary to manipulate the subscription and the underlying replication slot separately. Here are some scenarios:

- When creating a subscription, the replication slot already exists. In that case, the subscription can be created using the `create_slot = false` option to associate with the existing slot.
- When creating a subscription, the remote host is not reachable or in an unclear state. In that case, the subscription can be created using the `connect = false` option. The remote host will then not be contacted at all. This is what `pg_dump` uses. The remote replication slot will then have to be created manually before the subscription can be activated.
- When dropping a subscription, the replication slot should be kept. This could be useful when the subscriber database is being moved to a different host and will be activated from there. In that case, disassociate the slot from the subscription using `ALTER SUBSCRIPTION` before attempting to drop the subscription.
- When dropping a subscription, the remote host is not reachable. In that case, disassociate the slot from the subscription using `ALTER SUBSCRIPTION` before attempting to drop the subscription. If the remote database instance no longer exists, no further action is then necessary. If, however, the remote database instance is just unreachable, the replication slot (and any still remaining table synchronization slots) should then be dropped manually; otherwise it/they would continue to reserve WAL and might eventually cause the disk to fill up. Such cases should be carefully investigated.

29.2.2. Examples: Set Up Logical Replication

Create some test tables on the publisher.

```
test_pub=# CREATE TABLE t1(a int, b text, PRIMARY KEY(a));
CREATE TABLE
test_pub=# CREATE TABLE t2(c int, d text, PRIMARY KEY(c));
CREATE TABLE
test_pub=# CREATE TABLE t3(e int, f text, PRIMARY KEY(e));
CREATE TABLE
```

Create the same tables on the subscriber.


```
test_sub=# CREATE TABLE t1(a int, b text, PRIMARY KEY(a));
CREATE TABLE
test_sub=# CREATE TABLE t2(c int, d text, PRIMARY KEY(c));
CREATE TABLE
test_sub=# CREATE TABLE t3(e int, f text, PRIMARY KEY(e));
CREATE TABLE
```

Insert data to the tables at the publisher side.

```
test_pub=# INSERT INTO t1 VALUES (1, 'one'), (2, 'two'), (3, 'three');
INSERT 0 3
test_pub=# INSERT INTO t2 VALUES (1, 'A'), (2, 'B'), (3, 'C');
INSERT 0 3
test_pub=# INSERT INTO t3 VALUES (1, 'i'), (2, 'ii'), (3, 'iii');
INSERT 0 3
```

Create publications for the tables. The publications `pub2` and `pub3a` disallow some `publish` operations. The publication `pub3b` has a row filter (see [Section 29.4](#)).

```
test_pub=# CREATE PUBLICATION pub1 FOR TABLE t1;
CREATE PUBLICATION
test_pub=# CREATE PUBLICATION pub2 FOR TABLE t2 WITH (publish = 'truncate');
CREATE PUBLICATION
test_pub=# CREATE PUBLICATION pub3a FOR TABLE t3 WITH (publish = 'truncate');
CREATE PUBLICATION
test_pub=# CREATE PUBLICATION pub3b FOR TABLE t3 WHERE (e > 5);
CREATE PUBLICATION
```

Create subscriptions for the publications. The subscription `sub3` subscribes to both `pub3a` and `pub3b`. All subscriptions will copy initial data by default.

```
test_sub=# CREATE SUBSCRIPTION sub1
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=sub1'
test_sub=# PUBLICATION pub1;
CREATE SUBSCRIPTION
test_sub=# CREATE SUBSCRIPTION sub2
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=sub2'
test_sub=# PUBLICATION pub2;
CREATE SUBSCRIPTION
test_sub=# CREATE SUBSCRIPTION sub3
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=sub3'
test_sub=# PUBLICATION pub3a, pub3b;
CREATE SUBSCRIPTION
```

Observe that initial table data is copied, regardless of the `publish` operation of the publication.

```
test_sub=# SELECT * FROM t1;
 a |  b
---+-----
 1 | one
 2 | two
 3 | three
(3 rows)

test_sub=# SELECT * FROM t2;
 c |  d
---+---
 1 | A
 2 | B
 3 | C
(3 rows)
```

Furthermore, because the initial data copy ignores the `publish` operation, and because publication `pub3a` has no row filter, it means the copied table `t3` contains all rows even when they do not match the row filter of publication `pub3b`.

```
test_sub=# SELECT * FROM t3;
 e |  f
---+-----
 1 |  i
 2 | ii
 3 | iii
(3 rows)
```

Insert more data to the tables at the publisher side.

```
test_pub=# INSERT INTO t1 VALUES (4, 'four'), (5, 'five'), (6, 'six');
INSERT 0 3
test_pub=# INSERT INTO t2 VALUES (4, 'D'), (5, 'E'), (6, 'F');
INSERT 0 3
test_pub=# INSERT INTO t3 VALUES (4, 'iv'), (5, 'v'), (6, 'vi');
INSERT 0 3
```

Now the publisher side data looks like:

```
test_pub=# SELECT * FROM t1;
 a |  b
---+-----
 1 | one
 2 | two
 3 | three
 4 | four
 5 | five
 6 | six
(6 rows)
```

```
test_pub=# SELECT * FROM t2;
 c | d
---+---
 1 | A
 2 | B
 3 | C
 4 | D
 5 | E
 6 | F
(6 rows)
```

```
test_pub=# SELECT * FROM t3;
 e |  f
---+-----
 1 |  i
 2 | ii
 3 | iii
 4 | iv
 5 |  v
 6 | vi
(6 rows)
```

Observe that during normal replication the appropriate `publish` operations are used. This means publications `pub2` and `pub3a` will not replicate the `INSERT`. Also, publication `pub3b` will only replicate data that matches the row filter of `pub3b`. Now the subscriber side data looks like:

```
test_sub=# SELECT * FROM t1;
```

```

a |    b
---+-----
1 | one
2 | two
3 | three
4 | four
5 | five
6 | six
(6 rows)

```

```

test_sub=# SELECT * FROM t2;
c | d
---+---
1 | A
2 | B
3 | C
(3 rows)

```

```

test_sub=# SELECT * FROM t3;
e | f
---+-----
1 | i
2 | ii
3 | iii
6 | vi
(4 rows)

```

29.2.3. Examples: Deferred Replication Slot Creation

There are some cases (e.g. [Section 29.2.1](#)) where, if the remote replication slot was not created automatically, the user must create it manually before the subscription can be activated. The steps to create the slot and activate the subscription are shown in the following examples. These examples specify the standard logical decoding output plugin (`pgoutput`), which is what the built-in logical replication uses.

First, create a publication for the examples to use.

```

test_pub=# CREATE PUBLICATION pub1 FOR ALL TABLES;
CREATE PUBLICATION

```

Example 1: Where the subscription says `connect = false`

- Create the subscription.

```

test_sub=# CREATE SUBSCRIPTION sub1
test_sub=# CONNECTION 'host=localhost dbname=test_pub'
test_sub=# PUBLICATION pub1
test_sub=# WITH (connect=false);
WARNING:  subscription was created, but is not connected
HINT:     To initiate replication, you must manually create the replication slot,
         enable the subscription, and refresh the subscription.
CREATE SUBSCRIPTION

```

- On the publisher, manually create a slot. Because the name was not specified during `CREATE SUBSCRIPTION`, the name of the slot to create is same as the subscription name, e.g. "sub1".

```

test_pub=# SELECT * FROM pg_create_logical_replication_slot('sub1', 'pgoutput');
slot_name |    lsn
-----+-----
sub1      | 0/19404D0
(1 row)

```

- On the subscriber, complete the activation of the subscription. After this the tables of `pub1` will start replicating.

```
test_sub=# ALTER SUBSCRIPTION sub1 ENABLE;
ALTER SUBSCRIPTION
test_sub=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION;
ALTER SUBSCRIPTION
```

Example 2: Where the subscription says `connect = false`, but also specifies the `slot_name` option.

- Create the subscription.

```
test_sub=# CREATE SUBSCRIPTION sub1
test_sub=# CONNECTION 'host=localhost dbname=test_pub'
test_sub=# PUBLICATION pub1
test_sub=# WITH (connect=false, slot_name='myslot');
WARNING: subscription was created, but is not connected
HINT: To initiate replication, you must manually create the replication slot,
      enable the subscription, and refresh the subscription.
CREATE SUBSCRIPTION
```

- On the publisher, manually create a slot using the same name that was specified during `CREATE SUBSCRIPTION`, e.g. "myslot".

```
test_pub=# SELECT * FROM pg_create_logical_replication_slot('myslot', 'pgoutput');
 slot_name |      lsn
-----+-----
myslot    | 0/19059A0
(1 row)
```

- On the subscriber, the remaining subscription activation steps are the same as before.

```
test_sub=# ALTER SUBSCRIPTION sub1 ENABLE;
ALTER SUBSCRIPTION
test_sub=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION;
ALTER SUBSCRIPTION
```

Example 3: Where the subscription specifies `slot_name = NONE`

- Create the subscription. When `slot_name = NONE` then `enabled = false`, and `create_slot = false` are also needed.

```
test_sub=# CREATE SUBSCRIPTION sub1
test_sub=# CONNECTION 'host=localhost dbname=test_pub'
test_sub=# PUBLICATION pub1
test_sub=# WITH (slot_name=NONE, enabled=false, create_slot=false);
CREATE SUBSCRIPTION
```

- On the publisher, manually create a slot using any name, e.g. "myslot".

```
test_pub=# SELECT * FROM pg_create_logical_replication_slot('myslot', 'pgoutput');
 slot_name |      lsn
-----+-----
myslot    | 0/1905930
(1 row)
```

- On the subscriber, associate the subscription with the slot name just created.

```
test_sub=# ALTER SUBSCRIPTION sub1 SET (slot_name='myslot');
ALTER SUBSCRIPTION
```

- The remaining subscription activation steps are same as before.

```
test_sub=# ALTER SUBSCRIPTION sub1 ENABLE;
ALTER SUBSCRIPTION
test_sub=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION;
```

ALTER SUBSCRIPTION

29.3. Logical Replication Failover

To allow subscriber nodes to continue replicating data from the publisher node even when the publisher node goes down, there must be a physical standby corresponding to the publisher node. The logical slots on the primary server corresponding to the subscriptions can be synchronized to the standby server by specifying `failover = true` when creating subscriptions. See [Section 50.2.3](#) for details. Enabling the `failover` parameter ensures a seamless transition of those subscriptions after the standby is promoted. They can continue subscribing to publications on the new primary server.

Because the slot synchronization logic copies asynchronously, it is necessary to confirm that replication slots have been synced to the standby server before the failover happens. To ensure a successful failover, the standby server must be ahead of the subscriber. This can be achieved by configuring `synchronized_standby_slots`.

To confirm that the standby server is indeed ready for failover, follow these steps to verify that all necessary logical replication slots have been synchronized to the standby server:

1. On the subscriber node, use the following SQL to identify which replication slots should be synced to the standby that we plan to promote. This query will return the relevant replication slots associated with the failover-enabled subscriptions.

```
test_sub=# SELECT
           array_agg(quote_literal(s.subslotname)) AS slots
FROM   pg_subscription s
WHERE  s.subfailover AND
       s.subslotname IS NOT NULL;

 slots
-----
 {'sub1','sub2','sub3'}
(1 row)
```

2. On the subscriber node, use the following SQL to identify which table synchronization slots should be synced to the standby that we plan to promote. This query needs to be run on each database that includes the failover-enabled subscription(s). Note that the table sync slot should be synced to the standby server only if the table copy is finished (See [Section 54.57](#)). We don't need to ensure that the table sync slots are synced in other scenarios as they will either be dropped or re-created on the new primary server in those cases.

```
test_sub=# SELECT
           array_agg(quote_literal(slot_name)) AS slots
FROM   (
        SELECT CONCAT('pg_', srsubid, '_sync_', srrelid, '_',
ctl.system_identifier) AS slot_name
        FROM pg_control_system() ctl, pg_subscription_rel r, pg_subscription
s
        WHERE r.srsubstate = 'f' AND s.oid = r.srsubid AND s.subfailover
        );

 slots
-----
 {'pg_16394_sync_16385_7394666715149055164'}
(1 row)
```

3. Check that the logical replication slots identified above exist on the standby server and are ready for failover.

```
test_standby=# SELECT slot_name, (synced AND NOT temporary AND NOT conflicting) AS
failover_ready
FROM pg_replication_slots
```

```

WHERE slot_name IN
      ('sub1', 'sub2', 'sub3',
'pg_16394_sync_16385_7394666715149055164');
 slot_name | failover_ready
-----+-----
 sub1      | t
 sub2      | t
 sub3      | t
 pg_16394_sync_16385_7394666715149055164 | t
(4 rows)
    
```

If all the slots are present on the standby server and the result (`failover_ready`) of the above SQL query is true, then existing subscriptions can continue subscribing to publications now on the new primary server.

29.4. Row Filters

By default, all data from all published tables will be replicated to the appropriate subscribers. The replicated data can be reduced by using a *row filter*. A user might choose to use row filters for behavioral, security or performance reasons. If a published table sets a row filter, a row is replicated only if its data satisfies the row filter expression. This allows a set of tables to be partially replicated. The row filter is defined per table. Use a `WHERE` clause after the table name for each published table that requires data to be filtered out. The `WHERE` clause must be enclosed by parentheses. See [CREATE PUBLICATION](#) for details.

29.4.1. Row Filter Rules

Row filters are applied *before* publishing the changes. If the row filter evaluates to `false` or `NULL` then the row is not replicated. The `WHERE` clause expression is evaluated with the same role used for the replication connection (i.e. the role specified in the `CONNECTION` clause of the [CREATE SUBSCRIPTION](#)). Row filters have no effect for `TRUNCATE` command.

29.4.2. Expression Restrictions

The `WHERE` clause allows only simple expressions. It cannot contain user-defined functions, operators, types, and collations, system column references or non-immutable built-in functions.

If a publication publishes `UPDATE` or `DELETE` operations, the row filter `WHERE` clause must contain only columns that are covered by the replica identity (see [REPLICA IDENTITY](#)). If a publication publishes only `INSERT` operations, the row filter `WHERE` clause can use any column.

29.4.3. UPDATE Transformations

Whenever an `UPDATE` is processed, the row filter expression is evaluated for both the old and new row (i.e. using the data before and after the update). If both evaluations are `true`, it replicates the `UPDATE` change. If both evaluations are `false`, it doesn't replicate the change. If only one of the old/new rows matches the row filter expression, the `UPDATE` is transformed to `INSERT` or `DELETE`, to avoid any data inconsistency. The row on the subscriber should reflect what is defined by the row filter expression on the publisher.

If the old row satisfies the row filter expression (it was sent to the subscriber) but the new row doesn't, then, from a data consistency perspective the old row should be removed from the subscriber. So the `UPDATE` is transformed into a `DELETE`.

If the old row doesn't satisfy the row filter expression (it wasn't sent to the subscriber) but the new row does, then, from a data consistency perspective the new row should be added to the subscriber. So the `UPDATE` is transformed into an `INSERT`.

[Table 29.1](#) summarizes the applied transformations.

Table 29.1. UPDATE Transformation Summary

Old row	New row	Transformation
no match	no match	don't replicate
no match	match	INSERT
match	no match	DELETE
match	match	UPDATE

29.4.4. Partitioned Tables

If the publication contains a partitioned table, the publication parameter `publish_via_partition_root` determines which row filter is used. If `publish_via_partition_root` is `true`, the *root partitioned table's* row filter is used. Otherwise, if `publish_via_partition_root` is `false` (default), each *partition's* row filter is used.

29.4.5. Initial Data Synchronization

If the subscription requires copying pre-existing table data and a publication contains `WHERE` clauses, only data that satisfies the row filter expressions is copied to the subscriber.

If the subscription has several publications in which a table has been published with different `WHERE` clauses, rows that satisfy *any* of the expressions will be copied. See [Section 29.4.6](#) for details.

Warning

Because initial data synchronization does not take into account the `publish` parameter when copying existing table data, some rows may be copied that would not be replicated using DML. Refer to [Section 29.8.1](#), and see [Section 29.2.2](#) for examples.

Note

If the subscriber is in a release prior to 15, copy pre-existing data doesn't use row filters even if they are defined in the publication. This is because old releases can only copy the entire table data.

29.4.6. Combining Multiple Row Filters

If the subscription has several publications in which the same table has been published with different row filters (for the same `publish` operation), those expressions get ORed together, so that rows satisfying *any* of the expressions will be replicated. This means all the other row filters for the same table become redundant if:

- One of the publications has no row filter.
- One of the publications was created using `FOR ALL TABLES`. This clause does not allow row filters.
- One of the publications was created using `FOR TABLES IN SCHEMA` and the table belongs to the referred schema. This clause does not allow row filters.

29.4.7. Examples

Create some tables to be used in the following examples.

```
test_pub=# CREATE TABLE t1(a int, b int, c text, PRIMARY KEY(a,c));
CREATE TABLE
test_pub=# CREATE TABLE t2(d int, e int, f int, PRIMARY KEY(d));
```

```
CREATE TABLE
test_pub=# CREATE TABLE t3(g int, h int, i int, PRIMARY KEY(g));
CREATE TABLE
```

Create some publications. Publication p1 has one table (t1) and that table has a row filter. Publication p2 has two tables. Table t1 has no row filter, and table t2 has a row filter. Publication p3 has two tables, and both of them have a row filter.

```
test_pub=# CREATE PUBLICATION p1 FOR TABLE t1 WHERE (a > 5 AND c = 'NSW');
CREATE PUBLICATION
test_pub=# CREATE PUBLICATION p2 FOR TABLE t1, t2 WHERE (e = 99);
CREATE PUBLICATION
test_pub=# CREATE PUBLICATION p3 FOR TABLE t2 WHERE (d = 10), t3 WHERE (g = 10);
CREATE PUBLICATION
```

psql can be used to show the row filter expressions (if defined) for each publication.

```
test_pub=# \dRp+
                                Publication p1
  Owner   | All tables | Inserts | Updates | Deletes | Truncates | Via root
-----+-----+-----+-----+-----+-----+-----
postgres | f          | t       | t       | t       | t         | f
Tables:
    "public.t1" WHERE ((a > 5) AND (c = 'NSW'::text))

                                Publication p2
  Owner   | All tables | Inserts | Updates | Deletes | Truncates | Via root
-----+-----+-----+-----+-----+-----+-----
postgres | f          | t       | t       | t       | t         | f
Tables:
    "public.t1"
    "public.t2" WHERE (e = 99)

                                Publication p3
  Owner   | All tables | Inserts | Updates | Deletes | Truncates | Via root
-----+-----+-----+-----+-----+-----+-----
postgres | f          | t       | t       | t       | t         | f
Tables:
    "public.t2" WHERE (d = 10)
    "public.t3" WHERE (g = 10)
```

psql can be used to show the row filter expressions (if defined) for each table. See that table t1 is a member of two publications, but has a row filter only in p1. See that table t2 is a member of two publications, and has a different row filter in each of them.

```
test_pub=# \d t1
                                Table "public.t1"
 Column | Type   | Collation | Nullable | Default
-----+-----+-----+-----+-----
 a      | integer |           | not null |
 b      | integer |           |         |
 c      | text    |           | not null |
Indexes:
    "t1_pkey" PRIMARY KEY, btree (a, c)
Publications:
    "p1" WHERE ((a > 5) AND (c = 'NSW'::text))
    "p2"
```

```
test_pub=# \d t2
                                Table "public.t2"
```


Column	Type	Collation	Nullable	Default
d	integer		not null	
e	integer			
f	integer			

Indexes:

"t2_pkey" PRIMARY KEY, btree (d)

Publications:

"p2" WHERE (e = 99)

"p3" WHERE (d = 10)

test_pub=# \d t3

Table "public.t3"

Column	Type	Collation	Nullable	Default
g	integer		not null	
h	integer			
i	integer			

Indexes:

"t3_pkey" PRIMARY KEY, btree (g)

Publications:

"p3" WHERE (g = 10)

On the subscriber node, create a table t1 with the same definition as the one on the publisher, and also create the subscription s1 that subscribes to the publication p1.

```
test_sub=# CREATE TABLE t1(a int, b int, c text, PRIMARY KEY(a,c));
CREATE TABLE
test_sub=# CREATE SUBSCRIPTION s1
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=s1'
test_sub=# PUBLICATION p1;
CREATE SUBSCRIPTION
```

Insert some rows. Only the rows satisfying the t1 WHERE clause of publication p1 are replicated.

```
test_pub=# INSERT INTO t1 VALUES (2, 102, 'NSW');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (3, 103, 'QLD');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (4, 104, 'VIC');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (5, 105, 'ACT');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (6, 106, 'NSW');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (7, 107, 'NT');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (8, 108, 'QLD');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES (9, 109, 'NSW');
INSERT 0 1
```

```
test_pub=# SELECT * FROM t1;
```

a	b	c
2	102	NSW
3	103	QLD
4	104	VIC
5	105	ACT

```

6 | 106 | NSW
7 | 107 | NT
8 | 108 | QLD
9 | 109 | NSW
(8 rows)

```

```

test_sub=# SELECT * FROM t1;
 a |  b |  c
---+-----+-----
 6 | 106 | NSW
 9 | 109 | NSW
(2 rows)

```

Update some data, where the old and new row values both satisfy the `t1 WHERE` clause of publication `p1`. The `UPDATE` replicates the change as normal.

```

test_pub=# UPDATE t1 SET b = 999 WHERE a = 6;
UPDATE 1

```

```

test_pub=# SELECT * FROM t1;
 a |  b |  c
---+-----+-----
 2 | 102 | NSW
 3 | 103 | QLD
 4 | 104 | VIC
 5 | 105 | ACT
 7 | 107 | NT
 8 | 108 | QLD
 9 | 109 | NSW
 6 | 999 | NSW
(8 rows)

```

```

test_sub=# SELECT * FROM t1;
 a |  b |  c
---+-----+-----
 9 | 109 | NSW
 6 | 999 | NSW
(2 rows)

```

Update some data, where the old row values did not satisfy the `t1 WHERE` clause of publication `p1`, but the new row values do satisfy it. The `UPDATE` is transformed into an `INSERT` and the change is replicated. See the new row on the subscriber.

```

test_pub=# UPDATE t1 SET a = 555 WHERE a = 2;
UPDATE 1

```

```

test_pub=# SELECT * FROM t1;
 a |  b |  c
---+-----+-----
 3 | 103 | QLD
 4 | 104 | VIC
 5 | 105 | ACT
 7 | 107 | NT
 8 | 108 | QLD
 9 | 109 | NSW
 6 | 999 | NSW
555 | 102 | NSW
(8 rows)

```

```

test_sub=# SELECT * FROM t1;
 a |  b |  c

```

```
-----+-----+-----
  9 | 109 | NSW
  6 | 999 | NSW
555 | 102 | NSW
(3 rows)
```

Update some data, where the old row values satisfied the `t1 WHERE` clause of publication `p1`, but the new row values do not satisfy it. The `UPDATE` is transformed into a `DELETE` and the change is replicated. See that the row is removed from the subscriber.

```
test_pub=# UPDATE t1 SET c = 'VIC' WHERE a = 9;
UPDATE 1
```

```
test_pub=# SELECT * FROM t1;
 a | b | c
-----+-----+-----
  3 | 103 | QLD
  4 | 104 | VIC
  5 | 105 | ACT
  7 | 107 | NT
  8 | 108 | QLD
  6 | 999 | NSW
555 | 102 | NSW
  9 | 109 | VIC
(8 rows)
```

```
test_sub=# SELECT * FROM t1;
 a | b | c
-----+-----+-----
  6 | 999 | NSW
555 | 102 | NSW
(2 rows)
```

The following examples show how the publication parameter `publish_via_partition_root` determines whether the row filter of the parent or child table will be used in the case of partitioned tables.

Create a partitioned table on the publisher.

```
test_pub=# CREATE TABLE parent(a int PRIMARY KEY) PARTITION BY RANGE(a);
CREATE TABLE
test_pub=# CREATE TABLE child PARTITION OF parent DEFAULT;
CREATE TABLE
```

Create the same tables on the subscriber.

```
test_sub=# CREATE TABLE parent(a int PRIMARY KEY) PARTITION BY RANGE(a);
CREATE TABLE
test_sub=# CREATE TABLE child PARTITION OF parent DEFAULT;
CREATE TABLE
```

Create a publication `p4`, and then subscribe to it. The publication parameter `publish_via_partition_root` is set as `true`. There are row filters defined on both the partitioned table (`parent`), and on the partition (`child`).

```
test_pub=# CREATE PUBLICATION p4 FOR TABLE parent WHERE (a < 5), child WHERE (a >= 5)
test_pub=# WITH (publish_via_partition_root=true);
CREATE PUBLICATION

test_sub=# CREATE SUBSCRIPTION s4
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=s4'
test_sub=# PUBLICATION p4;
CREATE SUBSCRIPTION
```

Insert some values directly into the `parent` and `child` tables. They replicate using the row filter of `parent` (because `publish_via_partition_root` is `true`).

```
test_pub=# INSERT INTO parent VALUES (2), (4), (6);
INSERT 0 3
test_pub=# INSERT INTO child VALUES (3), (5), (7);
INSERT 0 3
```

```
test_pub=# SELECT * FROM parent ORDER BY a;
 a
---
 2
 3
 4
 5
 6
 7
(6 rows)
```

```
test_sub=# SELECT * FROM parent ORDER BY a;
 a
---
 2
 3
 4
(3 rows)
```

Repeat the same test, but with a different value for `publish_via_partition_root`. The publication parameter `publish_via_partition_root` is set as `false`. A row filter is defined on the partition (`child`).

```
test_pub=# DROP PUBLICATION p4;
DROP PUBLICATION
test_pub=# CREATE PUBLICATION p4 FOR TABLE parent, child WHERE (a >= 5)
test_pub=# WITH (publish_via_partition_root=false);
CREATE PUBLICATION

test_sub=# ALTER SUBSCRIPTION s4 REFRESH PUBLICATION;
ALTER SUBSCRIPTION
```

Do the inserts on the publisher same as before. They replicate using the row filter of `child` (because `publish_via_partition_root` is `false`).

```
test_pub=# TRUNCATE parent;
TRUNCATE TABLE
test_pub=# INSERT INTO parent VALUES (2), (4), (6);
INSERT 0 3
test_pub=# INSERT INTO child VALUES (3), (5), (7);
INSERT 0 3
```

```
test_pub=# SELECT * FROM parent ORDER BY a;
 a
---
 2
 3
 4
 5
 6
 7
(6 rows)
```

```
test_sub=# SELECT * FROM child ORDER BY a;
```

```
a
---
5
6
7
(3 rows)
```

29.5. Column Lists

Each publication can optionally specify which columns of each table are replicated to subscribers. The table on the subscriber side must have at least all the columns that are published. If no column list is specified, then all columns on the publisher are replicated. See [CREATE PUBLICATION](#) for details on the syntax.

The choice of columns can be based on behavioral or performance reasons. However, do not rely on this feature for security: a malicious subscriber is able to obtain data from columns that are not specifically published. If security is a consideration, protections can be applied at the publisher side.

If no column list is specified, any columns added to the table later are automatically replicated. This means that having a column list which names all columns is not the same as having no column list at all.

A column list can contain only simple column references. The order of columns in the list is not preserved.

Specifying a column list when the publication also publishes [FOR TABLES IN SCHEMA](#) is not supported.

For partitioned tables, the publication parameter [publish_via_partition_root](#) determines which column list is used. If [publish_via_partition_root](#) is `true`, the root partitioned table's column list is used. Otherwise, if [publish_via_partition_root](#) is `false` (the default), each partition's column list is used.

If a publication publishes `UPDATE` or `DELETE` operations, any column list must include the table's replica identity columns (see [REPLICA IDENTITY](#)). If a publication publishes only `INSERT` operations, then the column list may omit replica identity columns.

Column lists have no effect for the `TRUNCATE` command.

During initial data synchronization, only the published columns are copied. However, if the subscriber is from a release prior to 15, then all the columns in the table are copied during initial data synchronization, ignoring any column lists.

Warning: Combining Column Lists from Multiple Publications

There's currently no support for subscriptions comprising several publications where the same table has been published with different column lists. [CREATE SUBSCRIPTION](#) disallows creating such subscriptions, but it is still possible to get into that situation by adding or altering column lists on the publication side after a subscription has been created.

This means changing the column lists of tables on publications that are already subscribed could lead to errors being thrown on the subscriber side.

If a subscription is affected by this problem, the only way to resume replication is to adjust one of the column lists on the publication side so that they all match; and then either recreate the subscription, or use [ALTER SUBSCRIPTION ... DROP PUBLICATION](#) to remove one of the offending publications and add it again.

29.5.1. Examples

Create a table `t1` to be used in the following example.

```
test_pub=# CREATE TABLE t1(id int, a text, b text, c text, d text, e text, PRIMARY
KEY(id));
```

CREATE TABLE

Create a publication p1. A column list is defined for table t1 to reduce the number of columns that will be replicated. Notice that the order of column names in the column list does not matter.

```
test_pub=# CREATE PUBLICATION p1 FOR TABLE t1 (id, b, a, d);
CREATE PUBLICATION
```

psql can be used to show the column lists (if defined) for each publication.

```
test_pub=# \dRp+
                                Publication p1
 Owner   | All tables | Inserts | Updates | Deletes | Truncates | Via root
-----+-----+-----+-----+-----+-----+-----
postgres | f          | t       | t       | t       | t       | f
Tables:
    "public.t1" (id, a, b, d)
```

psql can be used to show the column lists (if defined) for each table.

```
test_pub=# \d t1
              Table "public.t1"
 Column | Type   | Collation | Nullable | Default
-----+-----+-----+-----+-----
 id     | integer |           | not null |
 a      | text    |           |          |
 b      | text    |           |          |
 c      | text    |           |          |
 d      | text    |           |          |
 e      | text    |           |          |
Indexes:
    "t1_pkey" PRIMARY KEY, btree (id)
Publications:
    "p1" (id, a, b, d)
```

On the subscriber node, create a table t1 which now only needs a subset of the columns that were on the publisher table t1, and also create the subscription s1 that subscribes to the publication p1.

```
test_sub=# CREATE TABLE t1(id int, b text, a text, d text, PRIMARY KEY(id));
CREATE TABLE
test_sub=# CREATE SUBSCRIPTION s1
test_sub=# CONNECTION 'host=localhost dbname=test_pub application_name=s1'
test_sub=# PUBLICATION p1;
CREATE SUBSCRIPTION
```

On the publisher node, insert some rows to table t1.

```
test_pub=# INSERT INTO t1 VALUES(1, 'a-1', 'b-1', 'c-1', 'd-1', 'e-1');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES(2, 'a-2', 'b-2', 'c-2', 'd-2', 'e-2');
INSERT 0 1
test_pub=# INSERT INTO t1 VALUES(3, 'a-3', 'b-3', 'c-3', 'd-3', 'e-3');
INSERT 0 1
test_pub=# SELECT * FROM t1 ORDER BY id;
 id | a  | b  | c  | d  | e
----+---+---+---+---+---
  1 | a-1 | b-1 | c-1 | d-1 | e-1
  2 | a-2 | b-2 | c-2 | d-2 | e-2
  3 | a-3 | b-3 | c-3 | d-3 | e-3
(3 rows)
```

Only data from the column list of publication p1 is replicated.

```
test_sub=# SELECT * FROM t1 ORDER BY id;
 id | b  | a  | d
-----+-----+-----+-----
  1 | b-1 | a-1 | d-1
  2 | b-2 | a-2 | d-2
  3 | b-3 | a-3 | d-3
(3 rows)
```

29.6. Conflicts

Logical replication behaves similarly to normal DML operations in that the data will be updated even if it was changed locally on the subscriber node. If incoming data violates any constraints the replication will stop. This is referred to as a *conflict*. When replicating `UPDATE` or `DELETE` operations, missing data will not produce a conflict and such operations will simply be skipped.

Logical replication operations are performed with the privileges of the role which owns the subscription. Permissions failures on target tables will cause replication conflicts, as will enabled [row-level security](#) on target tables that the subscription owner is subject to, without regard to whether any policy would ordinarily reject the `INSERT`, `UPDATE`, `DELETE` or `TRUNCATE` which is being replicated. This restriction on row-level security may be lifted in a future version of Postgres Pro.

A conflict will produce an error and will stop the replication; it must be resolved manually by the user. Details about the conflict can be found in the subscriber's server log.

The resolution can be done either by changing data or permissions on the subscriber so that it does not conflict with the incoming change or by skipping the transaction that conflicts with the existing data. When a conflict produces an error, the replication won't proceed, and the logical replication worker will emit the following kind of message to the subscriber's server log:

```
ERROR:  duplicate key value violates unique constraint "test_pkey"
DETAIL:  Key (c)=(1) already exists.
CONTEXT:  processing remote data for replication origin "pg_16395" during "INSERT" for
          replication target relation "public.test" in transaction 725 finished at 0/14C0378
```

The LSN of the transaction that contains the change violating the constraint and the replication origin name can be found from the server log (LSN 0/14C0378 and replication origin pg_16395 in the above case). The transaction that produced the conflict can be skipped by using `ALTER SUBSCRIPTION ... SKIP` with the finish LSN (i.e., LSN 0/14C0378). The finish LSN could be an LSN at which the transaction is committed or prepared on the publisher. Alternatively, the transaction can also be skipped by calling the `pg_replication_origin_advance()` function. Before using this function, the subscription needs to be disabled temporarily either by `ALTER SUBSCRIPTION ... DISABLE` or, the subscription can be used with the `disable_on_error` option. Then, you can use `pg_replication_origin_advance()` function with the `node_name` (i.e., pg_16395) and the next LSN of the finish LSN (i.e., 0/14C0379). The current position of origins can be seen in the `pg_replication_origin_status` system view. Please note that skipping the whole transaction includes skipping changes that might not violate any constraint. This can easily make the subscriber inconsistent.

When the `streaming` mode is `parallel`, the finish LSN of failed transactions may not be logged. In that case, it may be necessary to change the streaming mode to `on` or `off` and cause the same conflicts again so the finish LSN of the failed transaction will be written to the server log. For the usage of finish LSN, please refer to `ALTER SUBSCRIPTION ... SKIP`.

29.7. Restrictions

Logical replication currently has the following restrictions or missing functionality. These might be addressed in future releases.

- The database schema and DDL commands are not replicated. The initial schema can be copied by hand using `pg_dump --schema-only`. Subsequent schema changes would need to be kept in sync

manually. (Note, however, that there is no need for the schemas to be absolutely the same on both sides.) Logical replication is robust when schema definitions change in a live database: When the schema is changed on the publisher and replicated data starts arriving at the subscriber but does not fit into the table schema, replication will error until the schema is updated. In many cases, intermittent errors can be avoided by applying additive schema changes to the subscriber first.

- Sequence data is not replicated. The data in serial or identity columns backed by sequences will of course be replicated as part of the table, but the sequence itself would still show the start value on the subscriber. If the subscriber is used as a read-only database, then this should typically not be a problem. If, however, some kind of switchover or failover to the subscriber database is intended, then the sequences would need to be updated to the latest values, either by copying the current data from the publisher (perhaps using `pg_dump`) or by determining a sufficiently high value from the tables themselves.
- Replication of `TRUNCATE` commands is supported, but some care must be taken when truncating groups of tables connected by foreign keys. When replicating a truncate action, the subscriber will truncate the same group of tables that was truncated on the publisher, either explicitly specified or implicitly collected via `CASCADE`, minus tables that are not part of the subscription. This will work correctly if all affected tables are part of the same subscription. But if some tables to be truncated on the subscriber have foreign-key links to tables that are not part of the same (or any) subscription, then the application of the truncate action on the subscriber will fail.
- Large objects (see [Chapter 36](#)) are not replicated. There is no workaround for that, other than storing data in normal tables.
- Replication is only supported by tables, including partitioned tables. Attempts to replicate other types of relations, such as views, materialized views, or foreign tables, will result in an error.
- When replicating between partitioned tables, the actual replication originates, by default, from the leaf partitions on the publisher, so partitions on the publisher must also exist on the subscriber as valid target tables. (They could either be leaf partitions themselves, or they could be further sub-partitioned, or they could even be independent tables.) Publications can also specify that changes are to be replicated using the identity and schema of the partitioned root table instead of that of the individual leaf partitions in which the changes actually originate (see `publish_via_partition_root` parameter of `CREATE PUBLICATION`).
- When using `REPLICA IDENTITY FULL` on published tables, it is important to note that the `UPDATE` and `DELETE` operations cannot be applied to subscribers if the tables include attributes with datatypes (such as `point` or `box`) that do not have a default operator class for B-tree or Hash. However, this limitation can be overcome by ensuring that the table has a primary key or replica identity defined for it.

29.8. Architecture

Logical replication starts by copying a snapshot of the data on the publisher database. Once that is done, changes on the publisher are sent to the subscriber as they occur in real time. The subscriber applies data in the order in which commits were made on the publisher so that transactional consistency is guaranteed for the publications within any single subscription.

Logical replication is built with an architecture similar to physical streaming replication (see [Section 26.2.5](#)). It is implemented by `walsender` and `apply` processes. The `walsender` process starts logical decoding (described in [Chapter 50](#)) of the WAL and loads the standard logical decoding output plugin (`pgoutput`). The plugin transforms the changes read from WAL to the logical replication protocol (see [Section 56.5](#)) and filters the data according to the publication specification. The data is then continuously transferred using the streaming replication protocol to the `apply` worker, which maps the data to local tables and applies the individual changes as they are received, in correct transactional order.

The `apply` process on the subscriber database always runs with `session_replication_role` set to `replica`. This means that, by default, triggers and rules will not fire on a subscriber. Users can optionally choose to enable triggers and rules on a table using the `ALTER TABLE` command and the `ENABLE TRIGGER` and `ENABLE RULE` clauses.

The logical replication apply process currently only fires row triggers, not statement triggers. The initial table synchronization, however, is implemented like a `COPY` command and thus fires both row and statement triggers for `INSERT`.

29.8.1. Initial Snapshot

The initial data in existing subscribed tables are snapshotted and copied in a parallel instance of a special kind of apply process. This process will create its own replication slot and copy the existing data. As soon as the copy is finished the table contents will become visible to other backends. Once existing data is copied, the worker enters synchronization mode, which ensures that the table is brought up to a synchronized state with the main apply process by streaming any changes that happened during the initial data copy using standard logical replication. During this synchronization phase, the changes are applied and committed in the same order as they happened on the publisher. Once synchronization is done, control of the replication of the table is given back to the main apply process where replication continues as normal.

Note

The publication `publish` parameter only affects what DML operations will be replicated. The initial data synchronization does not take this parameter into account when copying the existing table data.

29.9. Monitoring

Because logical replication is based on a similar architecture as [physical streaming replication](#), the monitoring on a publication node is similar to monitoring of a physical replication primary (see [Section 26.2.5.2](#)).

The monitoring information about subscription is visible in `pg_stat_subscription`. This view contains one row for every subscription worker. A subscription can have zero or more active subscription workers depending on its state.

Normally, there is a single apply process running for an enabled subscription. A disabled subscription or a crashed subscription will have zero rows in this view. If the initial data synchronization of any table is in progress, there will be additional workers for the tables being synchronized. Moreover, if the `streaming` transaction is applied in parallel, there may be additional parallel apply workers.

29.10. Security

The role used for the replication connection must have the `REPLICATION` attribute (or be a superuser). If the role lacks `SUPERUSER` and `BYPASSRLS`, publisher row security policies can execute. If the role does not trust all table owners, include `options=-row_security=off` in the connection string; if a table owner then adds a row security policy, that setting will cause replication to halt rather than execute the policy. Access for the role must be configured in `pg_hba.conf` and it must have the `LOGIN` attribute.

In order to be able to copy the initial table data, the role used for the replication connection must have the `SELECT` privilege on a published table (or be a superuser).

To create a publication, the user must have the `CREATE` privilege in the database.

To add tables to a publication, the user must have ownership rights on the table. To add all tables in schema to a publication, the user must be a superuser. To create a publication that publishes all tables or all tables in schema automatically, the user must be a superuser.

There are currently no privileges on publications. Any subscription (that is able to connect) can access any publication. Thus, if you intend to hide some information from particular subscribers, such as by using row filters or column lists, or by not adding the whole table to the publication, be aware that other

publications in the same database could expose the same information. Publication privileges might be added to Postgres Pro in the future to allow for finer-grained access control.

To create a subscription, the user must have the privileges of the `pg_create_subscription` role, as well as `CREATE` privileges on the database.

The subscription apply process will, at a session level, run with the privileges of the subscription owner. However, when performing an insert, update, delete, or truncate operation on a particular table, it will switch roles to the table owner and perform the operation with the table owner's privileges. This means that the subscription owner needs to be able to `SET ROLE` to each role that owns a replicated table.

If the subscription has been configured with `run_as_owner = true`, then no user switching will occur. Instead, all operations will be performed with the permissions of the subscription owner. In this case, the subscription owner only needs privileges to `SELECT`, `INSERT`, `UPDATE`, and `DELETE` from the target table, and does not need privileges to `SET ROLE` to the table owner. However, this also means that any user who owns a table into which replication is happening can execute arbitrary code with the privileges of the subscription owner. For example, they could do this by simply attaching a trigger to one of the tables which they own. Because it is usually undesirable to allow one role to freely assume the privileges of another, this option should be avoided unless user security within the database is of no concern.

On the publisher, privileges are only checked once at the start of a replication connection and are not re-checked as each change record is read.

On the subscriber, the subscription owner's privileges are re-checked for each transaction when applied. If a worker is in the process of applying a transaction when the ownership of the subscription is changed by a concurrent transaction, the application of the current transaction will continue under the old owner's privileges.

29.11. Configuration Settings

Logical replication requires several configuration options to be set. Most options are relevant only on one side of the replication. However, `max_replication_slots` is used on both the publisher and the subscriber, but it has a different meaning for each.

29.11.1. Publishers

`wal_level` must be set to `logical`.

`max_replication_slots` must be set to at least the number of subscriptions expected to connect, plus some reserve for table synchronization.

`max_wal_senders` should be set to at least the same as `max_replication_slots`, plus the number of physical replicas that are connected at the same time.

Logical replication walsender is also affected by `wal_sender_timeout`.

29.11.2. Subscribers

`max_replication_slots` must be set to at least the number of subscriptions that will be added to the subscriber, plus some reserve for table synchronization.

`max_logical_replication_workers` must be set to at least the number of subscriptions (for leader apply workers), plus some reserve for the table synchronization workers and parallel apply workers.

`max_worker_processes` may need to be adjusted to accommodate for replication workers, at least $(\text{max_logical_replication_workers} + 1)$. Note, some extensions and parallel queries also take worker slots from `max_worker_processes`.

`max_sync_workers_per_subscription` controls the amount of parallelism of the initial data copy during the subscription initialization or when new tables are added.

`max_parallel_apply_workers_per_subscription` controls the amount of parallelism for streaming of in-progress transactions with subscription parameter `streaming = parallel`.

Logical replication workers are also affected by `wal_receiver_timeout`, `wal_receiver_status_interval` and `wal_retrieve_retry_interval`.

29.12. Quick Setup

First set the configuration options in `postgresql.conf`:

```
wal_level = logical
```

The other required settings have default values that are sufficient for a basic setup.

`pg_hba.conf` needs to be adjusted to allow replication (the values here depend on your actual network configuration and user you want to use for connecting):

```
host      all      repuser      0.0.0.0/0      md5
```

Then on the publisher database:

```
CREATE PUBLICATION mypub FOR TABLE users, departments;
```

And on the subscriber database:

```
CREATE SUBSCRIPTION mysub CONNECTION 'dbname=foo host=bar user=repuser' PUBLICATION  
mypub;
```

The above will start the replication process, which synchronizes the initial table contents of the tables `users` and `departments` and then starts replicating incremental changes to those tables.

Chapter 30. Just-in-Time Compilation (JIT)

This chapter explains what just-in-time compilation is, and how it can be configured in Postgres Pro.

30.1. What Is JIT compilation?

Just-in-Time (JIT) compilation is the process of turning some form of interpreted program evaluation into a native program, and doing so at run time. For example, instead of using general-purpose code that can evaluate arbitrary SQL expressions to evaluate a particular SQL predicate like `WHERE a.col = 3`, it is possible to generate a function that is specific to that expression and can be natively executed by the CPU, yielding a speedup.

Postgres Pro has builtin support to perform JIT compilation using [LLVM](#) when Postgres Pro is built with `--with-llvm`.

30.1.1. JIT Accelerated Operations

Currently Postgres Pro's JIT implementation has support for accelerating expression evaluation and tuple deforming. Several other operations could be accelerated in the future.

Expression evaluation is used to evaluate `WHERE` clauses, target lists, aggregates and projections. It can be accelerated by generating code specific to each case.

Tuple deforming is the process of transforming an on-disk tuple (see [Section 66.6.1](#)) into its in-memory representation. It can be accelerated by creating a function specific to the table layout and the number of columns to be extracted.

30.1.2. Inlining

Postgres Pro is very extensible and allows new data types, functions, operators and other database objects to be defined; see [Chapter 39](#). In fact the built-in objects are implemented using nearly the same mechanisms. This extensibility implies some overhead, for example due to function calls (see [Section 39.3](#)). To reduce that overhead, JIT compilation can inline the bodies of small functions into the expressions using them. That allows a significant percentage of the overhead to be optimized away.

30.1.3. Optimization

LLVM has support for optimizing generated code. Some of the optimizations are cheap enough to be performed whenever JIT is used, while others are only beneficial for longer-running queries. See <https://llvm.org/docs/Passes.html#transform-passes> for more details about optimizations.

30.2. When to JIT?

JIT compilation is beneficial primarily for long-running CPU-bound queries. Frequently these will be analytical queries. For short queries the added overhead of performing JIT compilation will often be higher than the time it can save.

To determine whether JIT compilation should be used, the total estimated cost of a query (see [Chapter 68](#) and [Section 19.7.2](#)) is used. The estimated cost of the query will be compared with the setting of `jit_above_cost`. If the cost is higher, JIT compilation will be performed. Two further decisions are then needed. Firstly, if the estimated cost is more than the setting of `jit_inline_above_cost`, short functions and operators used in the query will be inlined. Secondly, if the estimated cost is more than the setting of `jit_optimize_above_cost`, expensive optimizations are applied to improve the generated code. Each of these options increases the JIT compilation overhead, but can reduce query execution time considerably.

These cost-based decisions will be made at plan time, not execution time. This means that when prepared statements are in use, and a generic plan is used (see [PREPARE](#)), the values of the configuration parameters in effect at prepare time control the decisions, not the settings at execution time.

Note

If `jit` is set to `off`, or if no JIT implementation is available (for example because the server was compiled without `--with-llvm`), JIT will not be performed, even if it would be beneficial based on the above criteria. Setting `jit` to `off` has effects at both plan and execution time.

`EXPLAIN` can be used to see whether JIT is used or not. As an example, here is a query that is not using JIT:

```
=# EXPLAIN ANALYZE SELECT SUM(relpages) FROM pg_class;
                                         QUERY PLAN
-----
Aggregate  (cost=16.27..16.29 rows=1 width=8) (actual time=0.303..0.303 rows=1
 loops=1)
   -> Seq Scan on pg_class  (cost=0.00..15.42 rows=342 width=4) (actual
 time=0.017..0.111 rows=356 loops=1)
 Planning Time: 0.116 ms
 Execution Time: 0.365 ms
(4 rows)
```

Given the cost of the plan, it is entirely reasonable that no JIT was used; the cost of JIT would have been bigger than the potential savings. Adjusting the cost limits will lead to JIT use:

```
=# SET jit_above_cost = 10;
SET
=# EXPLAIN ANALYZE SELECT SUM(relpages) FROM pg_class;
                                         QUERY PLAN
-----
Aggregate  (cost=16.27..16.29 rows=1 width=8) (actual time=6.049..6.049 rows=1
 loops=1)
   -> Seq Scan on pg_class  (cost=0.00..15.42 rows=342 width=4) (actual
 time=0.019..0.052 rows=356 loops=1)
 Planning Time: 0.133 ms
 JIT:
   Functions: 3
   Options: Inlining false, Optimization false, Expressions true, Deforming true
   Timing: Generation 1.259 ms (Deform 0.000 ms), Inlining 0.000 ms, Optimization 0.797
 ms, Emission 5.048 ms, Total 7.104 ms
 Execution Time: 7.416 ms
```

As visible here, JIT was used, but inlining and expensive optimization were not. If `jit_inline_above_cost` or `jit_optimize_above_cost` were also lowered, that would change.

30.3. Configuration

The configuration variable `jit` determines whether JIT compilation is enabled or disabled. If it is enabled, the configuration variables `jit_above_cost`, `jit_inline_above_cost`, and `jit_optimize_above_cost` determine whether JIT compilation is performed for a query, and how much effort is spent doing so.

`jit_provider` determines which JIT implementation is used. It is rarely required to be changed. See [Section 30.4.2](#).

For development and debugging purposes a few additional configuration parameters exist, as described in [Section 19.20](#).

30.4. Extensibility

30.4.1. Inlining Support for Extensions

If you would like to build Postgres Pro extensions that support JIT inlining, make sure to set up your development environment as explained in [Section 17.1.5](#).

Postgres Pro's JIT implementation can inline the bodies of functions of types `C` and `internal`, as well as operators based on such functions. To do so for functions in extensions, the definitions of those functions need to be made available. When using [PGXS](#) to build an extension against a server that has been compiled with LLVM JIT support, the relevant files will be built and installed automatically.

The relevant files have to be installed into `$pkglibdir/bitcode/$extension/` and a summary of them into `$pkglibdir/bitcode/$extension.index.bc`, where `$pkglibdir` is the directory returned by `pg_config --pkglibdir` and `$extension` is the base name of the extension's shared library.

Note

For functions built into Postgres Pro itself, the bitcode is installed into `$pkglibdir/bitcode/postgres`.

30.4.2. Pluggable JIT Providers

Postgres Pro provides a JIT implementation based on LLVM. The interface to the JIT provider is pluggable and the provider can be changed without recompiling (although currently, the build process only provides inlining support data for LLVM). The active provider is chosen via the setting [jit_provider](#).

30.4.2.1. JIT Provider Interface

A JIT provider is loaded by dynamically loading the named shared library. The normal library search path is used to locate the library. To provide the required JIT provider callbacks and to indicate that the library is actually a JIT provider, it needs to provide a C function named `_PG_jit_provider_init`. This function is passed a struct that needs to be filled with the callback function pointers for individual actions:

```
struct JitProviderCallbacks
{
    JitProviderResetAfterErrorCB reset_after_error;
    JitProviderReleaseContextCB release_context;
    JitProviderCompileExprCB compile_expr;
};

extern void _PG_jit_provider_init(JitProviderCallbacks *cb);
```

Chapter 31. Enhanced Security

This chapter describes enhanced security mechanisms available in the certified edition of Postgres Pro Shardman. Moreover, Postgres Pro Shardman provides the following native security tools:

- The `pg_proaudit` extension enables detailed logging of various security events.
- Advanced authentication policies allow you to restrict database usage and enforce password requirements by assigning profiles to roles. For details on profiles, see `CREATE PROFILE`. The catalog `pg_profile` lists all profiles available for your cluster.

31.1. Memory Purge

Using `memory purge` configuration parameters, Postgres Pro Shardman can automatically replace data with zero bytes before it is deleted. This section describes how different types of data are handled. By default, all the memory purge parameters are switched on.

31.1.1. Deleting Files from External Memory

When you run some SQL commands that delete files, the corresponding disk space is returned to the operating system. Such commands include:

- Commands that delete database objects:
 - `DROP TABLE`
 - `DROP TEMPORARY TABLE`
 - `DROP MATERIALIZED VIEW`
 - `DROP INDEX`
 - `TRUNCATE`
 - `DROP DATABASE`
 - `DROP SCHEMA`
- Commands that re-create objects:
 - `VACUUM FULL`
 - `REINDEX`
 - `ALTER TABLE ADD COLUMN (with the default value)`
 - `ALTER TABLE ALTER COLUMN TYPE`

If the `wipe_file_on_delete` configuration parameter is switched on, the files to be deleted are first filled with zero bytes.

Note

The `ALTER TABLE DROP COLUMN` command does not re-create the file. All the data contained in the deleted column remains inside the pages, although you cannot access this data using SQL commands. To physically remove this data from the file, run `VACUUM FULL` after deleting the column.

31.1.2. Page Cleanup

In accordance with the Multiversion Concurrency Control (MVCC) model, when rows are deleted (with the `DELETE` command) the data stored in these rows is marked as deleted, but is not physically removed. In the case of row updates (with the `UPDATE` command), the old version of the row is deleted, and then a new version is inserted, so the previous value is not removed from the page either.

To ensure that MVCC mechanism is working correctly, Postgres Pro keeps the deleted data in the page while the saved version of the row can be accessed by at least one active snapshot. If the row is refer-

enced from indexes, index pages also keep references to the row versions that are deleted but not yet freed. When the row version is not referenced from any snapshot anymore, this version can be deleted with the `VACUUM` process. At the same time, all references to this row version are deleted from indexes. However, deleting does not mean physical removal: in normal operation, the corresponding space in the page is marked as free and can be used to store another row.

To prevent access to the deleted row versions, make sure that the `wipe_heaptuple_on_delete` parameter is switched on in the `postgresql.conf` configuration file. In this case, the `VACUUM` process not only marks the page spaces as free, but also fills them with zero bytes.

31.1.3. Clearing Random-Access Memory

While the server is running, random-access memory (RAM) is constantly allocated and released. Postgres Pro uses its own context-based system of memory allocation, in which memory is always allocated in one of the nested contexts. Deleting a context frees all the memory allocated in this context and the corresponding nested contexts. This approach helps to avoid memory leaks. Nevertheless, memory allocation and release are ultimately performed by the operating system. In normal operation, the released part of RAM is returned to the operating system and can be allocated to another process.

To remove all the data from the released part of RAM, make sure that the following parameters are switched on in the `postgresql.conf` configuration file:

- `wipe_memctx_on_free` — this parameter fills the released memory with zero bytes if the released memory belongs to a context.
- `wipe_mem_on_free` — this parameter fills the released memory with zero bytes if the released memory does not belong to any context. Although Postgres Pro always allocates memory within a particular context, you can still use this parameter to be on the safe side.

31.1.4. Cleaning up WAL Files

Write-Ahead Log (WAL) is a standard method for ensuring data integrity, which enables database backups, point-in-time recovery, and data replication between servers. WAL's central concept is that changes to data files (where tables and indexes reside) must be written only after those changes have been logged, that is, after log records describing the changes have been flushed to permanent storage. Thus, WAL can contain sensitive data that must be secured.

Postgres Pro server always stores a certain number of WAL segments in the `$PGDATA/pg_wal` directory. The minimum amount of WAL segments to store is defined by the `min_wal_size` parameter. In case of heavy load, the size of WAL segments can increase up to the `max_wal_size` value, or even exceed this value a little. As long as WAL disk usage stays above the `min_wal_size` threshold, old WAL segment files are deleted when they are released. Otherwise, WAL segments get overwritten.

To prevent unauthorized access to the released WAL segments, make sure that the `wipe_xlog_on_free` parameter is switched on in the `postgresql.conf` configuration file. In this case, the WAL segment will be filled with zero bytes before it is deleted or overwritten.

For details on WAL configuration and usage, see [Chapter 28](#).

31.2. Integrity Checks

31.2.1. Overview

Postgres Pro Shardman includes the `pg_integrity_check` utility that provides the following features for integrity checks:

- On-demand checksum calculation and validation
- Built-in checksum validation at the server startup

`pg_integrity_check` can control read-only files, additional files, and system catalog tables.

31.2.1.1. Read-Only Files

Controlled read-only files include executable files, libraries, and other files that must never be modified. Checksums for read-only files are managed by the `/opt/pgpro/sdm-17/share/security/system.conf` configuration file.

The `system.conf` file is included into the Postgres Pro Shardman distribution. Each Postgres Pro Shardman instance has a single `system.conf` file. Each line in `system.conf` corresponds to a single controlled object and includes two fields: the checksum, consisting of 40 hexadecimal digits, and a relative path to the controlled object. These fields are separated by three symbols: space, dash, space. Checksums are calculated and written into this file at the time of Postgres Pro Shardman installation. They control both contents and attributes of read-only files. For example, if access rights to the file have been modified, the checksum will change.

31.2.1.2. Additional Files

Additional files are controlled files that can be modified by a database administrator. Checksums for additional files are managed by configuration files stored in the `share/security` directory. Each cluster must have a separate configuration file for additional files. The naming convention is as follows: the path to the data directory of the cluster (PGDATA), with forward slashes replaced by underscores, followed by the `.user.conf` postfix. For example, for a cluster with `/var/lib/pgpro/sdm-17/data` data directory, the configuration file is called `_var_lib_pgpro_sdm-17_data.user.conf`.

Each line in the configuration file corresponds to a single controlled object and includes two fields: the checksum, consisting of 40 hexadecimal digits, and a relative path to the controlled object. The fields are separated by three symbols: space, dash, space.

Checksums control both contents and attributes of additional files. For example, if access rights to the file have been modified, the checksum will change.

To set up checksum validation for additional files, database administrator needs to do the following:

1. Create a configuration file for each cluster, following the naming convention specified above.
2. In the created configuration file, specify all the additional files to be controlled. As a checksum, any 40 hexadecimal digits can be entered, for example, zeros.
3. Run the following command to recalculate the checksums, specifying the path to the data directory of your cluster:

```
pg_integrity_check -u -o -D /var/lib/pgpro/sdm-17/data
```

Note

Alternatively, administrator can run the above command to generate a sample configuration file, and then edit this file as needed.

31.2.1.3. System Catalog Tables

System catalog tables represent a controlled data selection related to Postgres Pro Shardman instance. Checksums for system catalog tables are managed by respective configuration files for each database in a cluster. You can setup integrity checks of system catalog tables for several databases by providing separate configuration files with specific checksum values.

Each line in the configuration file corresponds to a single controlled object and includes two fields: the checksum, consisting of 40 hexadecimal digits, and a relative path to the controlled object. The fields are separated by three symbols: space, dash, space.

To set up checksum validation for data selection, database administrator needs to do the following for each database in a cluster:

1. Create a configuration file for the selected database. For example, *database-name-catalog.conf*.
2. In the created configuration file, specify SQL queries that return the controlled data. As a checksum, any 40 hexadecimal digits can be entered, for example, zeros.
3. Run the following command to recalculate checksums, specifying connection parameters for your database:

```
pg_integrity_check -c -o -C configuration-file-full-path -d postgres -h localhost -p
5432 -U postgres
```

To check system catalog tables, run `pg_integrity_check` for each database subject to integrity checks.

Note

Alternatively, administrator can run the above command to generate a sample configuration file, and then edit this file as needed.

31.2.2. Checking Integrity at Server Startup

Checksums for read-only files are always validated at the Postgres Pro Shardman server startup. If the server start is blocked because of a checksum mismatch, administrator must troubleshoot and resolve the issue to restart the server.

System catalog tables and additional files are not validated at the server startup. You can validate them by running `pg_integrity_check` manually once the server is started.

31.2.3. Scheduling Integrity Checks

If you are using a Linux-based system, you can schedule recurrent integrity checks using the `cron` daemon. To schedule integrity checks, modify the `/etc/crontab` file by adding the lines that define the frequency of `pg_integrity_check` utility launches. The `/etc/crontab` is a system file that contains all instructions for the `cron` daemon. To get a detailed description of the `crontab` file format for your Linux-based operating system, run the command:

```
man 5 crontab
```

Examples

The following example illustrates integrity checks on a Rosa SX operating system:

```
# Data directory of the cluster
PGDATA = /var/lib/pgpro/sdm-17/data
# pg_integrity_check log file
LOG = /opt/pgpro/sdm-17/share/security/log
# Validate read-only files every day at 00:05
5 0 * * *      root    /opt/pgpro/sdm-17/bin/pg_integrity_check -s >> $LOG
# Run integrity checks at 14:15 on the first day of the month
15 14 1 * *    root    /opt/pgpro/sdm-17/bin/pg_integrity_check -s >> $LOG
# Run integrity checks at 22.00 on weekdays
0 22 * * 1-5   root    /opt/pgpro/sdm-17/bin/pg_integrity_check -s >> $LOG
# Run integrity checks at 00:23, 2:23, 4:23 ..., every day
23 0-23/2 * * * root    /opt/pgpro/sdm-17/bin/pg_integrity_check -s >> $LOG
# Run integrity checks at 4:05 each Sunday
5 4 * * sun    root    /opt/pgpro/sdm-17/bin/pg_integrity_check -s >> $LOG
```

31.3. Separation of Duties between Privileged DBMS Users

The configuration of all PostgreSQL-based DBMS has the `superuser` role, which is needed at the initial database loading, but when later used for regular database operations, it carries information security

risks. As the `superuser` account has a wide range of system privileges, it can become a target for cyber-criminals and can cause adverse consequences, such as:

- Unauthorized access to sensitive data.
- Data leak.
- Dangerous changes to the DBMS configuration.
- DBMS failures.

Among the measures to resist such hazards, we can consider separation of duties, with reducing the number of operations that require superuser privileges. This is necessary to address risks both related to the trust in the superuser and to access of malicious users to the superuser account.

31.3.1. Overview

The Postgres Professional company has developed a model that separates superuser duties between two additional administrative roles — DBMS Administrator and Database Administrator (DB Administrator). Besides, mechanisms of protection against self-elevation of privileges and audit of all users' operations have been strengthened.

The DBMS Administrator is responsible for:

- The server management.
- Setup of the data replication and backups.
- Creation of databases.
- Setup of connections with the Postgres Pro database.
- Assignment of DB Administrators.

The DB Administrator is responsible for:

- Backup of the specific database.
- Creation of tables and other objects within the specific database.
- Creation of database users.
- Granting users with access rights.

Once most of the superuser's regular tasks are delegated to the new controlled administrators, the company can give up active use of this highly risky role. To do this, the infrastructure administrator must add rows that block superuser connections to the `pg_hba.conf` file and must also block updating this file by other administrators.

If superuser privileges must be temporarily returned to solve some rare complicated problem, this can be done in a single mode involving the infrastructure administrator. To switch the DBMS to a single mode, it should be stopped and then restarted as follows:

```
postgres --single -D DBMS_data_directory other_options DB_name
```

See [Section 17.2](#) for how to install additional supplied modules and extensions without superuser permissions.

The `pg_proaudit` extension, which enables logging various security events, is described in [Section F.6](#).

31.3.2. Creation of Additional Administrators

31.3.2.1. Creating the DBMS Administrator Role

This section describes creation of the `PGPRO_DBMS_ADMIN` role with the appropriate permissions.

The following script creates the `PGPRO_DBMS_ADMIN` role and the user with this role. It is run by the superuser `postgres`. Note that the script uses these [predefined roles](#):

- `pg_create_tablespace` allows executing the `CREATE TABLESPACE` command without superuser rights.
- `pg_manage_profiles` allows executing the `CREATE PROFILE`, `ALTER PROFILE`, and `DROP PROFILE` commands without superuser rights.

```
$ psql postgres
CREATE ROLE PGPRO_DBMS_ADMIN WITH CREATEDB CREATEROLE INHERIT LOGIN REPLICATION;
REVOKE ALL ON DATABASE postgres FROM PUBLIC;
GRANT CONNECT ON DATABASE postgres TO PGPRO_DBMS_ADMIN;
CREATE USER dbms_admin WITH CREATEDB CREATEROLE INHERIT LOGIN REPLICATION; -- Do not
forget to set initial password
GRANT PGPRO_DBMS_ADMIN TO dbms_admin;
GRANT pg_read_all_settings TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_read_all_stats TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_stat_scan_tables TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_monitor TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_signal_backend TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_checkpoint TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_create_tablespace TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
GRANT pg_manage_profiles TO PGPRO_DBMS_ADMIN WITH ADMIN OPTION;
EXIT;
```

The following commands grant permissions to manage the configuration and to perform logging and backup/restore operations to the `PGPRO_DBMS_ADMIN` role. It is run by the superuser `postgres`.

```
$ psql postgres
GRANT EXECUTE ON FUNCTION pg_reload_conf TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_rotate_logfile TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_create_restore_point TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_backup_start TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_backup_stop TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_switch_wal TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_promote TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_wal_replay_pause TO PGPRO_DBMS_ADMIN;
GRANT EXECUTE ON FUNCTION pg_wal_replay_resume TO PGPRO_DBMS_ADMIN;
EXIT;
```

31.3.2.2. Creating the DB Administrator Role

The `PGPRO_DB_DBNAME_ADMIN` role can include several users for one database. One user can also be included in several different `PGPRO_DB_DBNAME_ADMIN` roles if the user's job description includes management of several different databases.

This section describes creation of a database by example of the `DB1` database, its DB Administrator role and the `DB1_ADMIN` user granted appropriate rights. At this point, rights of the user with the `PGPRO_DBMS_ADMIN` role, which was created earlier, are sufficient for all the involved operations.

The following script creates the `DB1` database, the `PGPRO_DB_1_ADMIN` role and the user with this role. It is run by the DBMS Administrator:

```
$ psql postgres -U dbms_admin
SET ROLE PGPRO_DBMS_ADMIN;
CREATE ROLE PGPRO_DB_1_ADMIN WITH CREATEROLE INHERIT;
CREATE USER db1_admin WITH CREATEROLE INHERIT; -- Do not forget to set initial password
GRANT PGPRO_DB_1_ADMIN TO db1_admin;
GRANT PGPRO_DB_1_ADMIN TO PGPRO_DBMS_ADMIN;
GRANT pg_read_all_settings TO PGPRO_DB_1_ADMIN;
```

```
GRANT pg_read_all_stats TO PGPRO_DB_1_ADMIN;  
GRANT pg_stat_scan_tables TO PGPRO_DB_1_ADMIN;  
GRANT pg_monitor TO PGPRO_DB_1_ADMIN;  
CREATE DATABASE db1 OWNER PGPRO_DB_1_ADMIN;  
REVOKE CONNECT, TEMPORARY ON DATABASE db1 FROM PUBLIC;
```

31.3.3. Revoking Superuser Access

To block superuser connections, the infrastructure administrator adds these lines to `pg_hba.conf`:

TYPE	DATABASE	USER	ADDRESS	METHOD
local	all	postgres		reject
host	all	postgres	127.0.0.1/32	reject
host	all	postgres	:::1/128	reject

The root infrastructure administrator disallows the `postgres` user to change `pg_hba.conf`:

```
# chown root pg_hba.conf  
# chmod 640 pg_hba.conf
```

As a result, the `postgres` superuser cannot solely regain the possibility to connect to the DBMS:

```
# ls -lh pg_hba.conf  
-rw-r----- 1 root postgres ..... pg_hba.conf
```

Restart of the DBMS server is needed for the new restrictions to take effect:

```
# pg_ctl restart
```

31.4. Restricting DBMS Administrator's Data Access

Based on the [Separation of Duties between Privileged DBMS Users](#) technology, the Postgres Professional company has created a mechanism to block DBMS Administrator's and DB Administrator's data access, including viewing, but allow them to continue performing administration such as backup/restore, diagnostics, DB creation, deletion, and configuring.

31.4.1. Overview

Sensitive information is stored in tables, which are contained in DB schemas. By default, the DBMS Administrator and DB Administrators of appropriate databases can access this information, which makes a real issue. To protect data, Postgres Pro Shardman provides an enhanced security mechanism — the *Postgres Pro secured schema* and defines access rules for this schema.

This schema, described in [Section 5.10.3](#), will have a dedicated schema owner and security officer. They must be different DB users:

- The security officer (manager of access rights) will only be able to grant or revoke user access privileges for schema objects, but will be unable to work with schema objects: create, view, or modify them.
- The schema owner will be able to work with schema objects.

A superuser must create the owner and security officer of the secured schema. Hence, to create them and grant privileges, it is needed to switch the DBMS to a single mode and temporarily allow login of a superuser under the control of the trusted infrastructure administrator.

31.4.2. Creating Administrators and Managers of the DB Secured Schema

31.4.2.1. Granting a Temporary `postgres` Superuser Right for a DB Login from the Local Console

The infrastructure administrator `root` adds lines to `pg_hba.conf` that allow superuser connections from the local console:

TYPE	DATABASE	USER	ADDRESS	METHOD
local	all	postgres		peer
host	all	postgres	127.0.0.1/32	reject
host	all	postgres	:::1/128	reject

DBMS server restart is needed for the new restrictions to take effect:

```
# pg_ctl restart
```

31.4.2.2. Creating Managers of the Secured Schema

The process of creating a dedicated owner (manager) of the secured schema `DV_OWNER` and a security officer (manager of access rights) `DV_SEC_OFFICER` requires involvement of a superuser because if these accounts are created by the DBMS Administrator or DB Administrator, these administrators will still have the `ADMIN OPTION` right, which they even cannot revoke.

To address risks both related to the trust in administrators, who would retain the possibility to manage `DV_OWNER` and `DV_SEC_OFFICER`, and to access of malicious users to their accounts, the superuser performs the operations needed.

To simplify, this example assumes that the secured schema will be created in the `DB1` database.

```
$ psql postgres
CREATE USER DV_OWNER WITH LOGIN; -- Do not forget to set initial password
CREATE USER DV_SEC_OFFICER WITH LOGIN; -- Do not forget to set initial password
GRANT CONNECT ON DATABASE db1 TO DV_SEC_OFFICER;
GRANT CONNECT ON DATABASE db1 TO DV_OWNER;
EXIT;
```

31.4.2.3. Creating the Secured Schema

Creation of the secured schema is a critical activity, which must involve the superuser, who needs a temporary permission for the DB login.

For a role to be able to create an object that will be owned by a different role or to grant ownership of an existing object to a different role, the `SET ROLE` right is needed, otherwise the command `ALTER ... OWNER TO` or `CREATE DATABASE ... OWNER` will issue an error, as explained in [GRANT](#).

As the `DV_OWNER` is independent, access of the DBMS Administrator or DB Administrator to this role through `SET ROLE` should not be granted, and these administrators should not have ownership rights for SQL objects. Hence, creation of the secured schema by the superuser is the best way.

In the commands below, pay attention to the modified syntax of the `ALTER SCHEMA` command. The `SECURITY OFFICER/RESET SECURITY OFFICER` clause sets/resets the security officer, respectively. When you set the security officer, the schema becomes the `vault` schema, as described in [ALTER SCHEMA](#).

```
$ psql db1 -U postgres
CREATE SCHEMA vault;
GRANT ALL ON SCHEMA vault TO DV_OWNER;
GRANT USAGE ON SCHEMA vault TO DV_SEC_OFFICER;
ALTER SCHEMA vault OWNER TO DV_OWNER;
ALTER SCHEMA vault SECURITY OFFICER TO DV_SEC_OFFICER;
EXIT;
```

31.4.2.4. Protecting Secured Schema Users from Administrators

Users with access to the secured schema must be strongly protected. You cannot permit DBMS Administrator and DB Administrator have the [ADMIN OPTION](#) for secured schema users as in this case, these administrators may potentially get unauthorized access to the secured schema data.

To prevent from this vulnerability, the `ADMIN OPTION` for secured schema users must be revoked from other users by the [REVOKE ADMIN OPTION](#) command. Let's see how to do this by example of the `vault_user` user, which will be created by the `PGPRO_DB_1_ADMIN` role:

```
$ psql db1 -U db1_admin
CREATE USER vault_user; -- Do not forget to set passwords
GRANT CONNECT ON DATABASE db1 TO vault_user;
EXIT;
```

First, let's check who has ADMIN OPTION for vault_user:

```
$ psql db1 -U postgres
SELECT rn.rolname "role", mn.rolname member, gn.rolname grantor, am.admin_option
FROM pg_auth_members am
JOIN pg_roles rn ON rn.oid = am.roleid
JOIN pg_roles mn ON mn.oid = am.member
JOIN pg_roles gn ON gn.oid = am.grantor
WHERE rn.rolname = 'vault_user' AND am.admin_option = true
ORDER BY member;
```

In this example, the db1_admin role has these rights:

role	member	grantor	admin_option
vault_user	db1_admin	postgres	t

(1 row)

The postgres superuser will now revoke them from the db1_admin role:

```
$ psql db1 -U postgres
REVOKE ADMIN OPTION FOR vault_user FROM db1_admin CASCADE;
```

Finally, let's make sure that there are no extra users with the ADMIN OPTION for vault_user:

```
$ psql db1 -U postgres
SELECT rn.rolname "role", mn.rolname member, gn.rolname grantor, am.admin_option
FROM pg_auth_members am
JOIN pg_roles rn ON rn.oid = am.roleid
JOIN pg_roles mn ON mn.oid = am.member
JOIN pg_roles gn ON gn.oid = am.grantor
WHERE rn.rolname = 'vault_user' AND am.admin_option = true
ORDER BY member;
```

role	member	grantor	admin_option
------	--------	---------	--------------

(0 rows)

31.4.2.5. Revoking Superuser Access to the Secured Schema

To block superuser connections, the infrastructure administrator adds these lines to pg_hba.conf:

TYPE	DATABASE	USER	ADDRESS	METHOD
local	all	postgres		reject
host	all	postgres	127.0.0.1/32	reject
host	all	postgres	:::1/128	reject

Restart of the DBMS server is needed for the new restrictions to take effect:

```
# pg_ctl restart
```

31.4.3. Creating a Table in the Secured Schema

As the schema manager (owner) can only work with its objects, the DV_OWNER creates tables there:

```
$ psql db1 -U dv_owner
CREATE TABLE vault.vault_table (user_name TEXT, balance DECIMAL(10,2));
EXIT;
```

31.4.4. Moving a Table to the Secured Schema

Let's assume that the DB Administrator created a table:

```
$ psql db1 -U db1_admin
```

```
CREATE TABLE working_table(id INT, info TEXT);
```

This user can also fill this table with data:

```
INSERT INTO working_table(id, info) VALUES (1, 'Number One');
```

Later it was decided to move this table to the secured schema. To do this, the table ownership should be transferred to `DV_OWNER`, who will move it to the secured schema. However, PostgreSQL has a restricting requirement: the ownership of one's table can only be transferred to a user if there is a role that is member of both the old and new owning role:

```
ALTER TABLE working_table OWNER TO DV_OWNER;  
ERROR: must be able to SET ROLE "DV_OWNER"
```

Therefore, the superuser `postgres` must transfer the table ownership. To do this, the superuser `postgres` needs temporary DBMS access:

```
$ psql db1 -U postgres  
ALTER TABLE working_table OWNER TO DV_OWNER;
```

After execution of this command, DBMS access must be revoked from the superuser `postgres`.

`DV_OWNER` will move the table to the secured schema:

```
$ psql db1 -U dv_owner  
ALTER TABLE working_table SET SCHEMA vault;
```

31.4.5. Permitting Access to a Protected Table

As the schema security officer (manager of access rights) can only grant users with permissions to access schema objects, the `DV_SEC_OFFICER` manages the privileges. Only users that are explicitly specified will be able to access data. Other users, including DBMS Administrator and DB Administrator, will not have access to the data. Appropriate `GRANT` commands are as follows:

```
$ psql db1 -U dv_sec_officer  
GRANT USAGE ON SCHEMA vault TO vault_user;  
GRANT INSERT,UPDATE,DELETE,SELECT,TRUNCATE ON vault.vault_table TO vault_user;  
GRANT INSERT,UPDATE,DELETE,SELECT,TRUNCATE ON vault.working_table TO vault_user;  
EXIT;
```

31.4.6. Checking Access: Filling in a Protected Table

Connecting to the `db1` database as a user that was granted access to the secured schema, we can make sure that the user's privileges are sufficient to access data in the secured schema:

```
$ psql db1 -U vault_user  
INSERT INTO vault.vault_table( user_name, balance) VALUES( 'Ann', 200);  
INSERT INTO vault.vault_table( user_name, balance) VALUES( 'Bob', 3000);  
SELECT * FROM vault.vault_table;  
SELECT * FROM vault.working_table;  
EXIT;
```

31.4.7. Changes Made to a Backup Application

When restoring a secured schema, two different administrators are needed as one of them will be unable to restore access rights, while the other one will be unable to restore schema objects. Therefore, `pg_dump` now enables separate dumping of the schema and data and of access privileges for the secured schema. `--no-privileges` and `--privileges-only` options, respectively, enable this separate dumping.

Dump operations with `pg_dump` must be performed in the following order:

1. The DB Administrator dumps all data except the secured schema, using the `--exclude-schema` option.
2. The `DV_OWNER` dumps the secured schema data, using the `--schema` and `--no-privileges` options.
3. The `DV_OWNER` dumps the access privileges, using the `--schema` and `--privileges-only` options.

Restore operations using `pg_restore` must be performed in the following order:

1. The DB Administrator restores all data except the secured schema.
2. The `DV_OWNER` restores the secured schema data.
3. The `DV_SEC_OFFICER` restores the access privileges for the secured schema.

31.4.8. Changes Made to the Integrity-check Utility

The following changes have been made to `pg_integrity_check`:

- The new `-C` option allows you to specify the database whose system catalog must be validated. Tip: specify the one where the secured schema is created.
- The new `--syslog` option allows you to write checksum validation results into `syslog`.
- `pg_integrity_check` can now log information about both successful validation with the `Information` importance level and unsuccessful validation with the `Critical` importance level.

Chapter 32. Compressed File System (CFS)

This chapter describes the Compressed File System (CFS) feature, which enables page level compression in Postgres Pro Shardman.

32.1. Why database compression may be useful

Databases are used to store larger number of text and duplicated information. This is why compression of most of databases can be quite efficient and reduce used storage size 3..5 times. Postgres Pro Shardman performs compression of TOAST data, but small text fields which fits in the page are not compressed. Also not only heap pages can be compressed, indexes on text keys or indexes with larger number of duplicate values are also good candidates for compression.

Postgres Pro Shardman is working with disk data through buffer pool which accumulates most frequently used buffers. Interface between buffer manager and file system is the most natural place for performing compression. Buffers are stored on the disk in compressed form for reducing disk usage and minimizing amount of data to be read. And in-memory buffer pool contains uncompressed buffers, providing access to the records at the same speed as without compression. As far as modern server have large enough size of RAM, substantial part of the database can be cached in memory and accessed without any compression overhead penalty.

Except obvious advantage: saving disk space, compression can also improve system performance. There are two main reasons for it:

Reducing amount of disk IO

Compression helps to reduce size of data which should be written to the disk or read from it. Compression ratio 3 actually means that you need to read 3 times less data or same number of records can be fetched 3 times faster.

Improving locality

When modified buffers are flushed from buffer pool to the disk, they are written to the random locations on the disk. Postgres Pro Shardman cache replacement algorithm makes a decision about throwing away buffer from the pool based on its access frequency and ignoring its location on the disk. So two subsequently written buffers can be located in completely different parts of the disk. For HDD seek time is quite large - about 10msec, which corresponds to 100 random writes per second. And speed of sequential write can be about 100Mb/sec, which corresponds to 10000 buffers per second (100 times faster). For SSD gap between sequential and random write speed is smaller, but still sequential writers are more efficient. How it relates to data compression? Size of buffer in Postgres Pro Shardman is fixed (8kb by default). Size of compressed buffer depends on the content of the buffer. So updated buffer can not always fit in its old location on the disk. This is why we can not access pages directly by its address. Instead of it we have to use map which translates logical address of the page to its physical location on the disk. Definitely this extra level of indirection adds overhead. But in most cases this map can fit in memory, so page lookup is nothing more than just accessing array element. But presence of this map also have positive effect: we can now write updated pages sequentially, just updating their map entries. Postgres Pro Shardman is doing much to avoid "write storm" intensive flushing of data to the disk when buffer pool space is exhausted. Compression allows you to significantly reduce disk load.

32.2. How compression is integrated in Postgres Pro Shardman

To improve efficiency of disk IO, Postgres Pro Shardman is working with files through buffer manager, which pins in memory most frequently used pages. Each page is fixed size (8kb by default). But if we compress page, then its size will depend on its content. So updated page can require more (or less) space than original page. So we may not always perform in-place update of the page. Instead of it we have to locate new space for the page and somehow release old space. There are two main approaches to solving this problem:

Memory allocator

We should implement our own allocator of file space. Usually, to reduce fragmentation, fixed size block allocator is used. It means that we allocate place using some fixed quantum. For example if compressed page size is 932 bytes, then we will allocate 1024 bytes for it in the file.

Garbage collector

We can always allocate space for the pages sequentially at the end of the file and periodically do compactification (defragmentation) of the file, moving all used pages to the beginning of the file. Such garbage collection process can be performed in background. As it was explained in the previous section, sequential write of the flushed pages can significantly increase IO speed and some increase performance. This is why we have used this approach in CFS.

As far as page location is not fixed and page can be moved, we can not anymore access the page directly by its address and need to use extra level of indirection to map logical address of the page to it's physical location on the disk. It is done using mapping files. In most cases the mapping will be kept in memory (size of the map is 1000 times smaller size of the file) and address translation adds almost no overhead to page access time. But we need to maintain these extra files: flush them during checkpoint, remove when table is dropped, include them in backup and so on...

Postgres Pro Shardman stores relation in a set of files, size of each file is not exceeding 2Gb. Separate page map is constructed for each file. Garbage collection in CFS is done by several background workers. Number of this workers and pauses in their work can be configured by database administrator. These workers are splitting work based on inode hash, so them do not conflict with each other. Each file is proceeded separately. The file is blocked for access at the time of garbage collection but complete relation is not blocked. To ensure data consistency GC creates copies of original data and map files. Once them are flushed to the disk, new version of data file is atomically renamed to original file name. And then new page map data is copied to memory-mapped file and backup file for page map is removed. In case of recovery after crash we first inspect if there is backup of data file. If such file exists, then original file is not yet updated and we can safely remove backup files. If such file doesn't exist, then we check for presence of map file backup. If it is present, then defragmentation of this file was not completed because of crash and we complete this operation by copying map from backup file.

CFS supports zstd, zlib, and pglz compression libraries. On modern operating systems, lz4 compression is also supported. By default, zstd is used. For safety reasons, CFS checks that the compression algorithm used in a tablespace corresponds to this library.

32.3. Using Compression

Compression can only be enabled for separate tablespaces. To compress a tablespace, you should enable the `compression` option when creating this tablespace. For example:

```
postgres=# CREATE TABLESPACE tcfs LOCATION '/var/data/cfs' WITH (compression=true);
```

All tables created in this tablespace will be compressed using zstd, which is the default compression library.

Aside from the boolean value of the option, you can explicitly specify the library to use for compression. Possible values are `zstd`, `default` (the same as `zstd`), `pglz`, `zlib`, and `lz4`. For example, to use `zlib`, create the tablespace as follows:

```
postgres=# CREATE TABLESPACE tcfs1 LOCATION '/var/data/cfs1' WITH
(compression='zlib');
```

Once set, the tablespace compression option cannot be altered, so you cannot compress or decompress an already existing tablespace. While user relations are always compressed in CFS, system relations are not compressed if the numeric part of the filename, such as, for example, `pg_relation_filenode()`, is less than 16384 (see [Table 9.101](#) for more details).

If you would like to compress all tables created within the current session, you can make the compressed tablespace your default tablespace, as explained in [Section 22.7](#).

Note

For compressed tablespaces, `pg_checksums` will not verify checksums even if they are enabled.

To configure CFS, use configuration parameters listed in [Section 19.16](#). By default, CFS launches one background worker performing garbage collection. Garbage collector traverses tablespace directory, locating map files in it and checking percent of garbage in this file. When ratio of used and allocated spaces exceeds `cfs_gc_threshold` threshold, this file is defragmented. The file is locked at the period of defragmentation, preventing any access to this part of relation. To avoid getting stuck during any file defragmentation failure, CFS waits `cfs_gc_respond_time` seconds before the file is released from the lock. If it is not, a warning is written to the log. When defragmentation is completed, garbage collection waits `cfs_gc_delay` milliseconds and continues directory traversal. After the end of traversal, GC waits `cfs_gc_period` milliseconds and starts new GC iteration. If there are more than one GC workers, then they split work based on hash of file inode.

Several functions allow you to calculate the disk space used by different objects in CFS. For example:

- `pg_relation_size()` can compute the disk space used by the relation in CFS.
- `pg_total_relation_size()` computes the total disk space used by the specified table in CFS, including all indexes and TOAST data.
- `pg_indexes_size()` computes the total disk space used by indexes attached to the specified table in CFS.
- `pg_database_size()` computes the disk space used by the specified database in CFS.

See [Section 9.28.7](#) for more details.

CFS provides several functions to manually control CFS garbage collection and get information on CFS state and activity. For the full list of functions, see [Section 9.28.11](#).

To initiate garbage collection manually, use the `cfs_start_gc(n_workers)` function. This function returns number of workers which are actually started. Please notice that if `cfs_gc_workers` parameter is non zero, then GC is performed in background and `cfs_start_gc` function does nothing and returns 0.

Like the automatic garbage collection, the `cfs_start_gc(n_workers)` function only processes relations if the percent of garbage blocks in this relation exceeds the `cfs_gc_threshold` value. To defragment a relation with a smaller percent of garbage, you can temporarily set this parameter to a smaller value for your current session before calling this function.

It is possible to estimate effect of table compression using `cfs_estimate(relation)` function. This function takes first ten blocks of relation, tries to compress them, and returns average compress ratio. So if returned value is 7.8 then compressed table occupies about eight time less space than original table.

Function `cfs_compression_ratio(relation)` allows you to check how precise was the estimation of `cfs_estimate(relation)` function. It returns real compression ration for all segments of the compressed relation. Compression ration is total sum of virtual size of all relation segments (number of blocks multiplied by 8kb) divided by sum of physical size of the segment files.

As it was mentioned before, CFS always appends updated blocks to the end of the compressed file. So physical size of the file can be greater than used size in this file. I.e. CFS file is fragmented and defragmentation is periodically performed by CFS garbage collector. `cfs_fragmentation(relation)` functions returns average fragmentation of relation files. It is calculated as sum of physical sizes of the files minus sum of used size of the files divided by sum of physical sizes of the files.

To perform defragmentation for a particular compressed relation, use the `cfs_gc_relation(relation)` function. It returns the number of processed segments of the relation. Just like garbage collection per-

formed in the background, this function only processes segments in which the percent of garbage blocks exceeds the `cfs_gc_threshold` value.

There are several functions allowing to monitor garbage collection activity: `cfs_gc_activity_scanned_files` returns number of files scanned by GC, `cfs_gc_activity_processed_files` returns number of file compacted by GC, `cfs_gc_activity_processed_pages` returns number of pages transferred by GC during files defragmentation, `cfs_gc_activity_processed_bytes` returns total size of transferred pages. All these functions calculate their values since system start.

Chapter 33. Cluster Configuration for Distributed System

Postgres Pro Shardman provides a wide range of high-availability services and tools to cover the needs of most high-load systems.

33.1. Cluster Initialization Settings Related to Access Management

When a Postgres Pro Shardman cluster is initialized, security-related settings are taken from the initialization file. You can change them later, but do this with care and remember that in most cases, the change will require a DBMS restart.

A Postgres Pro Shardman cluster has two special users: *administrative* and *replication*. Postgres Pro Shardman manages controlled DBMS instances with administrative users. BiHA needs replication users for replications between controlled DBMS instances.

Security-related settings from the initialization file specify:

- Authentication methods for administrative and replication users — `PgSuAuthMethod`, `PgReplAuthMethod`
- Usernames for administrative and replication users — `PgSuUsername`, `PgReplUsername`
- Passwords for administrative and replication users — `PgSuPassword`, `PgReplPassword`
- `pg_hba.conf` rules used by DBMS instances — `StolonSpec.pgHBA`

See [Section 19.21.2](#) for detailed descriptions of these settings.

To change security-related user settings, perform these steps:

1. Check that the user that you want to specify in `PgReplUsername` / `PgSuUsername` exists with `REPLICATION/ SUPERUSER` privileges on all replication groups in the cluster and his password matches the new `PgReplPassword` / `PgSuPassword` setting.
2. If this is true, create dump of the `shardman/cluster0/data/cluster` etcd key (here and further the name of the Postgres Pro Shardman cluster is assumed to be `cluster0`). For example:

```
$ etcdctl --endpoints etcdserver:2379 get --print-value-only shardman/cluster0/data/cluster | jq . > clusterdata.json
```

This example creates the dump of the `data/cluster` key for the Postgres Pro Shardman cluster with the `cluster0` name from the etcd server `etcdserver` listening on port 2379, formats the dump with `jq` and saves to the `clusterdata.json` file.

3. Edit the dump as necessary and store it back in etcd:

```
$ cat clusterdata.json | etcdctl --endpoints etcdserver:2379 put shardman/cluster0/data/cluster
```

Modifying these settings will lead to a DBMS restart.

Unlike the above settings, the `StolonSpec.pgHBA` setting can be changed online. To do this, perform these steps:

1. Extract the `StolonSpec` definition from `shardman/cluster0/data/cluster`, save to some file, modify as necessary and update cluster settings with the `shardmanctl config update` command:

```
$ etcdctl --endpoints etcdserver:2379 get --print-value-only shardman/cluster0/data/cluster | jq .Spec.StolonSpec . > shardspec.json
```

2. Edit `shardspec.json` and replace the `StolonSpec.pgHBA` definition with the appropriate one, for example:

3. Apply the edited shardspec.json file:

The Postgres Pro Shardman architecture allows you to scale out your cluster without any downtime. This section describes how you can add more nodes to your Postgres Pro Shardman cluster in order to improve query performance/scalability. If a Postgres Pro Shardman cluster does not meet your performance expectations or storage capacity, you can add new nodes to the cluster.

In the *manual-topology* mode, to add a primary to a cluster, use the `shardmanctl nodes add` command, which adds the list of nodes to the cluster as primaries with a separate replication group for each primary. Create a cluster with three primary nodes and manual topology (`PlacementPolicy=manual` in `sdmspec.json`):

To view the topology of a cluster, use the `shardmanctl cluster topology` command:

The command output is as follows:

895

Cluster Configuration for Distributed System

PRIMARY			
		n2	5432
== REPLICATION GROUP clover-3-n3, RGID - 3 ==			
		HOST	PORT
PRIMARY			
		n3	5432
STATUS			

Add n4, n5, n6 nodes as replicas using the `shardmanctl shard add` command:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://
etcd2:2379,http://etcd3:2379 shard --shard clover-1-n1 add -n n4
$ shardmanctl --store-endpoints http://etcd1:2379,http://
etcd2:2379,http://etcd3:2379 shard --shard clover-2-n2 add -n n5
$ shardmanctl --store-endpoints http://etcd1:2379,http://
etcd2:2379,http://etcd3:2379 shard --shard clover-3-n3 add -n n6
```

In manual-topology mode, one node can be added to more than one replication group.

As a result, we get the following cluster configuration:

== REPLICATION GROUP clover-1-n1, RGID - 1 ==			
		HOST	PORT
		n1	5432
		n4	5432
PRIMARY			
STANDBY			
== REPLICATION GROUP clover-2-n2, RGID - 2 ==			
		HOST	PORT
STATUS			

Cluster Configuration for Distributed System

	n2	5432	PRIMARY
	n5	5432	STANDBY
== REPLICATION GROUP clover-3-n3, RGID - 3 ==			
	HOST	PORT	STATUS
	n3	5432	PRIMARY
	n6	5432	STANDBY

To remove a replica, just run the `shardmanctl shard rm` command. For example:

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://
etcd2:2379,http://etcd3:2379 shard --shard clover-1-n1 rm -n n4
```

To remove the master, first run the `shardmanctl shard switch` command to switch the master to the replica; then delete the old master.

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://
etcd2:2379,http://etcd3:2379 shard --shard clover-1-n1 switch --new-primary n4
```

33.3. Rebalancing the Data

33.3.1. Automatically Rebalancing the Data

Automatic rebalancing is used as the default mode. A rebalance process starts automatically after adding nodes (by default if `--no-rebalance` is not set) or before deleting a node. Rebalance can also be started manually. The essence of the rebalancing process is to evenly distribute partitions for each sharded table between replication groups.

The rebalancing process for each sharded table iteratively determines the replication group with the maximum and minimum number of partitions and creates a task to move one partition to the replication group with the minimum number of partitions. This process is repeated while $\text{max} - \text{min} > 1$. To move partitions, we use logical replication. Partitions of colocated tables are moved together with partitions of the sharded tables to which they refer.

It is important to remember that `max_logical_replication_workers` should be rather high since the rebalance process uses up to $\text{max}(\text{max_replication_slots}, \text{max_logical_replication_workers}, \text{max_worker_processes}, \text{max_wal_senders})/3$ concurrent threads. In practice, you can use `max_log-`

`ical_replication_workers = Repfactor + 3 * task_num` (`task_num` is the number of parallel rebalance tasks).

To rebalance sharded tables in the `cluster0` cluster manually, run the command (where `etcd1`, `etcd2`, `etcd3` are `etcd` cluster nodes):

```
$ shardmanctl --store-endpoints http://etcd1:2379,http://  
etcd2:2379,http://etcd3:2379 rebalance
```

If the process ends with an error, then you need to call the `shardmanctl cleanup` command with the `--after-rebalance` option.

33.3.2. Manually Rebalancing the Data

There are times when you need to place partitions of sharded tables in a specific way across the cluster nodes. To solve this problem, Postgres Pro Shardman supports the manual data rebalancing mode.

How it works:

1. Get a list of sharded tables using the `shardmanctl tables sharded list` command. As a result, we get an answer similar to the following:

```
$ shardmanctl shardmanctl tables sharded list
```

Sharded tables:

```
public.doc  
public.resolution  
public.users
```

2. Request information about the selected sharded tables. Example:

```
$ shardmanctl shardmanctl tables sharded info -t public.users
```

Table `public.users`

Partitions:

Partition	RgID	Shard	Master
0	1	clover-1-shrn1	shrn1:5432
1	2	clover-2-shrn2	shrn2:5432
2	3	clover-3-shrn3	shrn3:5432
3	1	clover-1-shrn1	shrn1:5432
4	2	clover-2-shrn2	shrn2:5432
5	3	clover-3-shrn3	shrn3:5432
6	1	clover-1-shrn1	shrn1:5432
7	2	clover-2-shrn2	shrn2:5432
8	3	clover-3-shrn3	shrn3:5432
9	1	clover-1-shrn1	shrn1:5432
10	2	clover-2-shrn2	shrn2:5432
11	3	clover-3-shrn3	shrn3:5432
12	1	clover-1-shrn1	shrn1:5432
13	2	clover-2-shrn2	shrn2:5432
14	3	clover-3-shrn3	shrn3:5432
15	1	clover-1-shrn1	shrn1:5432
16	2	clover-2-shrn2	shrn2:5432
17	3	clover-3-shrn3	shrn3:5432
18	1	clover-1-shrn1	shrn1:5432
19	2	clover-2-shrn2	shrn2:5432

Cluster Configuration for Distributed System

20	3	clover-3-shrn3	shrn3:5432
21	1	clover-1-shrn1	shrn1:5432
22	2	clover-2-shrn2	shrn2:5432
23	3	clover-3-shrn3	shrn3:5432

3. Move a partition to a new shard, as shown below:

```
$ shardmanctl --log-level debug tables sharded partmove -t public.users --partnum 1
--shard clover-1-shrn1
```

```
2023-07-26T06:00:36.900Z      DEBUG    cmd/common.go:105      Waiting for metadata
lock...
2023-07-26T06:00:36.936Z      DEBUG    rebalance/service.go:256      take
extension lock
2023-07-26T06:00:36.938Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=3
2023-07-26T06:00:36.938Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=2
2023-07-26T06:00:36.938Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=1
2023-07-26T06:00:36.951Z      DEBUG    broadcaster/worker.go:51      repgroup 3
connect established
2023-07-26T06:00:36.951Z      DEBUG    broadcaster/worker.go:51      repgroup 2
connect established
2023-07-26T06:00:36.952Z      DEBUG    broadcaster/worker.go:51      repgroup 1
connect established
2023-07-26T06:00:36.952Z      DEBUG    extension/lock.go:35      Waiting for
extension lock...
2023-07-26T06:00:36.976Z      INFO     rebalance/service.go:276      Performing
move partition...
2023-07-26T06:00:36.977Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=3
2023-07-26T06:00:36.978Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=2
2023-07-26T06:00:36.978Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=1
2023-07-26T06:00:36.987Z      DEBUG    broadcaster/worker.go:51      repgroup 1
connect established
2023-07-26T06:00:36.989Z      DEBUG    broadcaster/worker.go:51      repgroup 2
connect established
2023-07-26T06:00:36.992Z      DEBUG    broadcaster/worker.go:51      repgroup 3
connect established
2023-07-26T06:00:36.992Z      DEBUG    rebalance/service.go:71      Performing cleanup
after possible rebalance operation failure
2023-07-26T06:00:37.077Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=3
2023-07-26T06:00:37.077Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=1
2023-07-26T06:00:37.077Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=2
2023-07-26T06:00:37.082Z      DEBUG    rebalance/service.go:422      Rebalance
will run 1 tasks
2023-07-26T06:00:37.095Z      DEBUG    rebalance/service.go:452      Guessing
that rebalance() can use 3 workers
2023-07-26T06:00:37.096Z      DEBUG    rebalance/job.go:352      state: Idle
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move"}
```

Cluster Configuration for Distributed System

```
2023-07-26T06:00:37.111Z      DEBUG    rebalance/job.go:352      state:
ConnEstablished {"worker_id": 1, "table": "users", "partition num": 1, "source
rgid": 2, "dest rgid": 1, "kind": "move"}
2023-07-26T06:00:37.171Z      DEBUG    rebalance/job.go:352      state: WaitInitCopy
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move"}
2023-07-26T06:00:38.073Z      DEBUG    rebalance/job.go:347      current state
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move", "state": "WaitInitialCatchup"}
2023-07-26T06:00:38.073Z      DEBUG    rebalance/job.go:352      state:
WaitInitialCatchup {"worker_id": 1, "table": "users", "partition num": 1,
"source rgid": 2, "dest rgid": 1, "kind": "move"}
2023-07-26T06:00:38.084Z      DEBUG    rebalance/job.go:347      current state
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move", "state": "WaitFullSync"}
2023-07-26T06:00:38.084Z      DEBUG    rebalance/job.go:352      state: WaitFullSync
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move"}
2023-07-26T06:00:38.108Z      DEBUG    rebalance/job.go:347      current state
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move", "state": "Committing"}
2023-07-26T06:00:38.108Z      DEBUG    rebalance/job.go:352      state: Committing
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move"}
2023-07-26T06:00:38.254Z      DEBUG    rebalance/job.go:352      state: Complete
{"worker_id": 1, "table": "users", "partition num": 1, "source rgid": 2, "dest
rgid": 1, "kind": "move"}
2023-07-26T06:00:38.258Z      DEBUG    rebalance/service.go:583      Produce and
process tasks on destination replication groups...
2023-07-26T06:00:38.258Z      DEBUG    rebalance/service.go:594      Produce and
process tasks on source replication groups...
2023-07-26T06:00:38.258Z      DEBUG    rebalance/service.go:606      wait all
tasks finish
2023-07-26T06:00:38.258Z      DEBUG    rebalance/service.go:531      Analyzing
table public.users in rg 1 {"table": "public.users", "rgid": 1, "action":
"analyze"}
2023-07-26T06:00:38.573Z      DEBUG    rebalance/service.go:531      Analyzing
table public.users in rg 2 {"table": "public.users", "rgid": 2, "action":
"analyze"}
2023-07-26T06:00:38.833Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=1
2023-07-26T06:00:38.833Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=2
2023-07-26T06:00:38.833Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=3
```

In this example, partition number 1 of the `public.users` table will be moved to the `clover-1-shrn1` shard.

After manually moving a partition of a sharded table and for all tables collocated with it, automatic data rebalancing for these tables will be disabled.

To get the list of tables with disabled automatic rebalancing, call the `shardmanctl tables sharded norebalance` command. Example:

```
$ shardmanctl tables sharded norebalance
```

public.users

To enable automatic data rebalancing for a selected sharded table, call the `shardmanctl tables sharded rebalance` command, as shown in the example below:

\$ shardmanctl tables sharded rebalance -t public.users

```

2023-07-26T07:07:00.657Z      DEBUG    cmd/common.go:105      Waiting for metadata
lock...
2023-07-26T07:07:00.687Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=1
2023-07-26T07:07:00.687Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=2
2023-07-26T07:07:00.687Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=3
2023-07-26T07:07:00.697Z      DEBUG    broadcaster/worker.go:51      repgroup 1
connect established
2023-07-26T07:07:00.698Z      DEBUG    broadcaster/worker.go:51      repgroup 2
connect established
2023-07-26T07:07:00.698Z      DEBUG    broadcaster/worker.go:51      repgroup 3
connect established
2023-07-26T07:07:00.698Z      DEBUG    extension/lock.go:35      Waiting for extension
lock...
2023-07-26T07:07:00.719Z      DEBUG    rebalance/service.go:381      Planned moving
pnum 21 for table users from rg 1 to rg 2
2023-07-26T07:07:00.719Z      INFO     rebalance/service.go:244      Performing
rebalance...
2023-07-26T07:07:00.720Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=1
2023-07-26T07:07:00.720Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=2
2023-07-26T07:07:00.720Z      DEBUG    broadcaster/worker.go:33      start
broadcaster worker for repgroup id=3
2023-07-26T07:07:00.732Z      DEBUG    broadcaster/worker.go:51      repgroup 3
connect established
2023-07-26T07:07:00.732Z      DEBUG    broadcaster/worker.go:51      repgroup 1
connect established
2023-07-26T07:07:00.734Z      DEBUG    broadcaster/worker.go:51      repgroup 2
connect established
2023-07-26T07:07:00.734Z      DEBUG    rebalance/service.go:71      Performing cleanup
after possible rebalance operation failure
2023-07-26T07:07:00.791Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=1
2023-07-26T07:07:00.791Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=2
2023-07-26T07:07:00.791Z      DEBUG    broadcaster/worker.go:75      finish
broadcaster worker for repgroup id=3
2023-07-26T07:07:00.795Z      DEBUG    rebalance/service.go:422      Rebalance will
run 1 tasks
2023-07-26T07:07:00.809Z      DEBUG    rebalance/service.go:452      Guessing that
rebalance() can use 3 workers
2023-07-26T07:07:00.809Z      DEBUG    rebalance/job.go:352      state: Idle
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}
2023-07-26T07:07:00.823Z      DEBUG    rebalance/job.go:352      state: ConnsEstablished
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}

```

```

2023-07-26T07:07:00.880Z          DEBUG    rebalance/job.go:352    state: WaitInitCopy
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}
2023-07-26T07:07:01.886Z          DEBUG    rebalance/job.go:347    current state
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move", "state": "WaitInitialCatchup"}
2023-07-26T07:07:01.886Z          DEBUG    rebalance/job.go:352    state:
WaitInitialCatchup      {"worker_id": 1, "table": "users", "partition num": 21,
"source rgid": 1, "dest rgid": 2, "kind": "move"}
2023-07-26T07:07:01.904Z          DEBUG    rebalance/job.go:347    current state
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move", "state": "WaitFullSync"}
2023-07-26T07:07:01.905Z          DEBUG    rebalance/job.go:352    state: WaitFullSync
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}
2023-07-26T07:07:01.932Z          DEBUG    rebalance/job.go:347    current state
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move", "state": "Committing"}
2023-07-26T07:07:01.932Z          DEBUG    rebalance/job.go:352    state: Committing
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}
2023-07-26T07:07:02.057Z          DEBUG    rebalance/job.go:352    state: Complete
{"worker_id": 1, "table": "users", "partition num": 21, "source rgid": 1, "dest rgid":
2, "kind": "move"}
2023-07-26T07:07:02.060Z          DEBUG    rebalance/service.go:583    Produce and
process tasks on destination replication groups...
2023-07-26T07:07:02.060Z          DEBUG    rebalance/service.go:594    Produce and
process tasks on source replication groups...
2023-07-26T07:07:02.060Z          DEBUG    rebalance/service.go:531    Analyzing table
public.users in rg 2      {"table": "public.users", "rgid": 2, "action": "analyze"}
2023-07-26T07:07:02.060Z          DEBUG    rebalance/service.go:606    wait all tasks
finish
2023-07-26T07:07:02.321Z          DEBUG    rebalance/service.go:531    Analyzing table
public.users in rg 1      {"table": "public.users", "rgid": 1, "action": "analyze"}
2023-07-26T07:07:02.587Z          DEBUG    broadcaster/worker.go:75    finish
broadcaster worker for repgroup id=3
2023-07-26T07:07:02.587Z          DEBUG    broadcaster/worker.go:75    finish
broadcaster worker for repgroup id=2
2023-07-26T07:07:02.587Z          DEBUG    broadcaster/worker.go:75    finish
broadcaster worker for repgroup id=1

```

To enable automatic data rebalancing for all sharded tables, run the `shardmanctl rebalance` command with the `--force` option.

```
$ shardmanctl rebalance --force
```

33.4. Fault Tolerance and High Availability

Postgres Pro Shardman provides out-of-the-box fault tolerance. The `shardmand` daemon monitors the cluster configuration and manages BiHA clusters, which are used to guarantee high availability of all shards and fault tolerance. The common Postgres Pro Shardman configuration (`shardmand`, BiHA clusters) is stored in an etcd cluster.

In the case of automatic failover, BiHA will take care of automatically changing slave to master and failed master to standby. Information about the current states of PostgreSQL instances is periodically stored in etcd and used by cluster software.

If necessary, you can switch to a new primary manually by running the `shardmanctl shard switch` command.

You can specify some high-availability options to define cluster behavior in a fault state: `masterDemotionEnabled`, `masterDemotionTimeout`, `minSyncMonitorEnabled` and `minSyncMonitorUnhealthyTimeout`.

33.5. Cluster Services

The Postgres Pro Shardman cluster configuration is stored in etcd. Postgres Pro Shardman cluster services are organized as systemd services. The Postgres Pro Shardman configuration daemon `shardmand` monitors the cluster configuration and manages PostgreSQL instances. Each node has one `shardmand` service, whose typical name is `shardmand@CLUSTER_NAME.service`. Here `CLUSTER_NAME` is the Postgres Pro Shardman cluster name, `cluster0` by default.

Each `shardmand` includes several integrated `keeper` threads.

Each registered DBMS instance has an associated `keeper` thread that directly manages this PostgreSQL instance. The `keeper` starts, stops, initializes and update PostgreSQL instances according to the desired cluster state. Postgres Pro Shardman only uses synchronous replicas (otherwise, there is a chance to lose data when a node fails).

`shardmand` is a systemd unit, its logs are written to `journald`. You can use `journalctl` to examine it. For example, to get all logs since `2023-05-09 10:00` for the `shardmand` service of the `cluster0` cluster, you can use the following command:

```
$ journalctl -u shardmand@cluster0 --since '2023-05-09 10:00'
```

To control the log verbosity for all Postgres Pro Shardman services, set `SDM_LOG_LEVEL` in the `shardmand` configuration file.

33.6. Configuring Secure Communications with etcd

This section describes how to configure secure communications between the etcd store and Postgres Pro Shardman services and tools.

`etcd` is a critical component for a Postgres Pro Shardman cluster. If an intruder gets access to the `etcd` store, it gains full control over the whole cluster, including DBMS access with DBA privileges. To protect your cluster, it is recommended that you configure TLS authentication between `etcd` daemons and Postgres Pro Shardman services.

To this end, you can use HTTPS transport with certificates signed by your local certification authority (CA) to encrypt traffic between the `etcd` cluster and Postgres Pro Shardman services and restrict `etcd` access. To do this, perform the steps described in the next sections.

33.6.1. Generating SSL Certificates

To generate SSL certificates, perform the following steps:

1. If the CA does not exist, generate a self-signed root certificate. Generate all certificates on one trusted host. Here certificates that expire in 10000 days are generated (you can choose a more suitable interval):

```
# openssl genrsa -out rootCA.key 4096
# openssl req -x509 -new -key rootCA.key -days 10000 -out rootCA.crt
```

2. Prepare the following `openssl` configuration file for each `etcd` host:

```
[ req ]
default_bits          = 4096
distinguished_name    = req_distinguished_name
req_extensions        = req_ext
[ req_distinguished_name ]
countryName            = Country Name (2 letter code)
```

```
stateOrProvinceName      = State or Province Name (full name)
localityName             = Locality Name (eg, city)
organizationName         = Organization Name (eg, company)
commonName               = Common Name (e.g. server FQDN or YOUR name)
[ req_ext ]
subjectAltName = @alt_names
[alt_names]
DNS.1      = n1
IP.1       = 192.168.1.1
IP.2       = 127.0.0.1
```

Under `[alt_names]`, specify alternative subject names for the etcd host. These names must include the etcd server hostname, IP address and local IP. Including the local IP is convenient rather than required.

Save the file. For example, the names of configuration files for nodes `n1` — `n3` can be `n1.san.conf` — `n3.san.conf`.

3. Using the configuration files prepared, generate private keys and certificate requests for etcd hosts:

```
# openssl genrsa -out n1.etcd.key 4096
# openssl req -config n1.san.conf -new -key n1.etcd.key -out n1.etcd.csr -subj "/
C=RU/ST=Moscow Region/L=Moscow/O=Test/CN=n1"
```

Here `" /C=RU/ST=Moscow Region/L=Moscow/O=Test/CN=n1"` means that the certificate request is generated with the country name `RU`, state `Moscow Region`, locality `Moscow`, organization `Test` and common name `n1`. The common name must match the DNS name of your etcd server.

4. Sign the certification request:

```
# openssl x509 -extfile n1.san.conf -extensions req_ext -req -in n1.etcd.csr -CA
rootCA.crt -CAkey rootCA.key -CAcreateserial -out n1.etcd.crt -days 10000
```

5. Check the certificates to ensure they contain correct X509v3 Subject Alternative Name fields. The fields must contain the list of DNS names and IP addresses that you added to the openssl configuration file:

```
# openssl x509 -in n1.etcd.crt -noout -text
```

6. Generate client certificates for Postgres Pro Shardman services and client tools. These certificates do not need to contain the `subjectAltName` header, and CN is not important in the example below. It generates one common certificate-key pair for services and one — for tools:

```
# openssl x509 genrsa -out shardman_services.key 4096
# openssl req -new -key shardman_services.key -out shardman_services.csr -subj "/
C=RU/ST=Moscow Region/L=Moscow/O=Test/CN=shardman_services"
# openssl x509 -req -in shardman_services.csr -CA rootCA.crt -CAkey rootCA.key -
CAcreateserial -out shardman_services.crt -days 10000
# openssl x509 genrsa -out shardman_tools.key 4096
# openssl req -new -key shardman_tools.key -out shardman_tools.csr -subj "/C=RU/
ST=Moscow Region/L=Moscow/O=Test/CN=shardman_tools"
# openssl x509 -req -in shardman_tools.csr -CA rootCA.crt -CAkey rootCA.key -
CAcreateserial -out shardman_tools.crt -days 10000
```

33.6.2. Configuring etcd and shardmand Services

Now configure services (`etcd` and `shardmand`) to use the generated certificates. To do this, perform the following steps:

1. On each etcd node, put `rootCA.crt`, `nX.etcd.crt` and `nX.etcd.key` in the location accessible to the etcd daemon (for example, create `/etc/etcd` directory and put files there). Ensure that the `nX.etcd.key` file is only accessible to the etcd daemon user.
2. Specify the following configuration for etcd daemons in `/etc/default/etcd`:


```
# unqualified first name
ETCD_NAME=n1
# where we actually listen for peers
ETCD_LISTEN_PEER_URLS=https://0.0.0.0:2380
# where we actually listen for clients
ETCD_LISTEN_CLIENT_URLS=https://0.0.0.0:2379
# advertise where this machine is listening for clients
ETCD_ADVERTISE_CLIENT_URLS=https://n1:2379

# --initial flags are used during bootstrapping and ignored afterwards, so it is
# ok to specify them always
# advertise where this machine is listening for peer
ETCD_INITIAL_ADVERTISE_PEER_URLS=https://n1:2380
ETCD_INITIAL_CLUSTER_TOKEN=etcd-cluster
# ansible_nodename is fqdn
ETCD_INITIAL_CLUSTER=n1=https://n1:2380,n2=https://n2:2380,n3=https://n3:2380
ETCD_INITIAL_CLUSTER_STATE=new

ETCD_DATA_DIR=/var/lib/etcd/default/member
ETCD_AUTO_COMPACTION_RETENTION=1

ETCD_KEY_FILE=/etc/etcd/n1.etcd.key
ETCD_CERT_FILE=/etc/etcd/n1.etcd.crt
ETCD_TRUSTED_CA_FILE=/etc/etcd/rootCA.crt
ETCD_CLIENT_CERT_AUTH=true

ETCD_PEER_CERT_FILE=/etc/etcd/n1.etcd.crt
ETCD_PEER_KEY_FILE=/etc/etcd/n1.etcd.key
ETCD_PEER_TRUSTED_CA_FILE=/etc/etcd/rootCA.crt
ETCD_PEER_CLIENT_CERT_AUTH=true
```

Replace n1 here with the appropriate node name.

3. Restart etcd services on all etcd cluster nodes:

```
# systemctl restart etcd
```

4. To check the new configuration, use the following command:

```
# etcdctl --endpoints=https://n1:2379,https://n2:2379,https://n3:2379 --cacert /  
etc/etcd/rootCA.crt --cert /etc/etcd/n1.etcd.crt --key /etc/etcd/n1.etcd.key member  
list -w table
```

```
+-----+-----+-----+-----+-----+
+-----+
|          ID          | STATUS | NAME |   PEER ADDRS   |   CLIENT ADDRS   | IS LEARNER |
|          |
+-----+-----+-----+-----+-----+
+-----+
| 66ebe06d7302c3f0 | started |  n2  | https://n2:2380 | https://n2:2379 |      false |
|          |
| b1080bf5ff059980 | started |  n1  | https://n1:2380 | https://n1:2379 |      false |
|          |
| d98323257249fefb | started |  n3  | https://n3:2380 | https://n3:2379 |      false |
|          |
+-----+-----+-----+-----+-----+
+-----+
```

5. On each Postgres Pro Shardman cluster node, put rootCA.crt, shardman_services.crt and shardman_services.key in a location accessible to the postgres user (for example, create the /etc/

shardman directory and put files there). Ensure that the `shardman_services.key` file is only accessible to the postgres user.

6. Edit the shardmand configuration file `/etc/shardman/shardmand-cluster0.env` as follows:

```
SDM_STORE_ENDPOINTS=https://n1:2379,https://n2:2379,https://n3:2379
SDM_STORE_CA_FILE=/etc/shardman/rootCA.crt
SDM_STORE_CERT_FILE=/etc/shardman/shardman_services.crt
SDM_STORE_KEY=/etc/shardman/shardman_services.key
```

7. Restart `shardmand@cluster0` services on all Shardman nodes:

```
# systemctl restart shardmand@cluster0
```

33.6.3. Using Postgres Pro Shardman Tools

Before using Postgres Pro Shardman tools, copy `rootCA.crt`, `shardman_tools.crt` and `shardman_tools.key` to some location on the Postgres Pro Shardman management node where they are accessible to the management user. Here, any node with installed Postgres Pro Shardman utilities that is used to manage the Postgres Pro Shardman cluster is meant by *management node*. This can also be one of the Shardman cluster nodes (or all of them). By *management user*, a user is meant who runs `shardmanctl` tool. It is assumed that the certificates and key are located in the `/etc/shardman` directory.

When using Postgres Pro Shardman tools, be sure to add `--store-ca-file`, `--store-cert-file` and `--store-key` options to `shardmanctl` command. For example, the following command gets the cluster status:

```
# shardmanctl --store-ca-file /etc/shardman/rootCA.crt --store-cert-file /etc/shardman/
shardman_tools.crt --store-key /etc/shardman/shardman_tools.key --store-endpoints
https://n1:2379,https://n2:2379,https://n3:2379 status
```

Chapter 34. Troubleshooting

Postgres Pro offers the ability to dump the state of a backend process, which can be useful for diagnostic and debugging purposes, by enabling the [crash_info](#) configuration parameter. To set it up, follow these steps:

- Create a directory on each cluster node that the Postgres Pro Shardman operating system user has access to (usually, it is `postgres`). Error reports will be sent to this directory.

```
install -d -o postgres -g postgres -m 700 /var/lib/postgresql/crashinfo
```

- Set the `crash_info_location` value.

Note

This will cause the DBMS to restart.

```
shardmanctl --store-endpoints http://etcdserver:2379 set -y crash_info_location=/var/lib/postgresql/crashinfo
```

- To make sure the changes are applied, send a signal that will cause the backend failure and a core dump creation, along with the instance restart.

Note

Do it in your test environment only.

Connect to your DBMS and find out PID of the backend associated with the current session:

```
postgres=# select pg_backend_pid();
pg_backend_pid
-----
23770
```

Then send the SIGSEGV signal to the process with the received PID:

```
kill -11 23770
```

This will result in this backend crash, and a log file with the time, backtrace and cause of an error will be written to `/var/lib/postgresql/crashinfo`:

```
# Signal
Program received signal: 11 (SIGSEGV)
Signal      UTC date time: 25.10.2024 08:37:02

# Program
                                pid: 23770
                                ppid: 17506
    program_invocation_name: postgres postgres postgres 10.42.42.10(34202) idle
program_invocation_short_name: tgres 10.42.42.10(34202) idle
                                exe_path: /opt/pgpro/sdm-17/bin/postgres
                                exe: postgres

# Backtrace
1  postgres + 0x5b55c0          0x55c5ba8459b7  0x00007ffcbef19070
   bt_crash_handler + 0x3f7
2  libc.so.6 + 0x4251f          0x7f01c2caa520  0x00007ffcbef19140  __sigaction +
   0x50
```

```
unknown ./signal/../sysdeps/unix/sysv/linux/x86_64/libc_sigaction.c:0
3  libc.so.6 + 0x125f80          0x7f01c2d8df9a  0x00007ffcbef195b8  epoll_wait +
  0x1a
epoll_wait ./sysdeps/unix/sysv/linux/epoll_wait.c:30
4  postgres + 0x433870          0x55c5ba6c39bb  0x00007ffcbef195c0
  WaitEventSetWait + 0x14b
5  postgres + 0x320de0          0x55c5ba5b0e74  0x00007ffcbef19650  secure_read +
  0x94
6  postgres + 0x327d20          0x55c5ba5b7dae  0x00007ffcbef196a0  pq_recvbuf +
  0x8e
7  postgres + 0x328980          0x55c5ba5b8995  0x00007ffcbef196c0  pq_getbyte +
  0x15
8  postgres + 0x457da0          0x55c5ba6e909c  0x00007ffcbef196d0  PostgresMain +
  0x12fc
9  postgres + 0x3ce210          0x55c5ba65ef86  0x00007ffcbef19a60  ServerLoop +
  0xd76
10 postgres + 0x3cf240          0x55c5ba65fe18  0x00007ffcbef1a040  PostmasterMain
  + 0xbd8
11 postgres + 0x14ecc0          0x55c5ba3df182  0x00007ffcbef1a0c0  main + 0x4c2
12 libc.so.6 + 0x29d10          0x7f01c2c91d90  0x00007ffcbef1a0f0
  __libc_init_first + 0x90
__libc_start_call_main ./sysdeps/nptl/libc_start_call_main.h:58
13 libc.so.6 + 0x29dc0          0x7f01c2c91e40  0x00007ffcbef1a190
  __libc_start_main + 0x80
call_init ./csu/libc-start.c:128
__libc_start_main_impl ./csu/libc-start.c:379
14 postgres + 0x14f200          0x55c5ba3df225  0x00007ffcbef1a1e0  _start + 0x25
```

The dump state file can be generated in one of the following ways:

- By sending the signal 40 (also known as the diagnostic dump signal):

```
kill -40 backend_pid
```

- Using the `pg_diagdump()` function:

```
SELECT pg_diagdump(backend_pid);
```

Here `backend_pid` is the process ID of the backend process to dump.

As a result, Postgres Pro will write the state dump to a file in the `$PGDATA/crash_info` directory by default or in the directory specified in the [crash_info_location](#) configuration parameter. The file will be named following this pattern: `crash_file_id_pidpid.state`. You can set the data sources to provide data for a crash dump in the [crash_info_dump](#) configuration parameter.

The below example shows how to generate and inspect the state dump file for the backend with PID 23111:

```
-- Generate the state dump file
SELECT pg_diagdump(23111);

-- Inspect crash_info directory and its contents
SELECT pg_ls_dir('crash_info');

-- Read the contents of the state dump file
SELECT pg_read_file('crash_info/crash_1722943138419104_pid23111.state');
```

Part IV. Client Interfaces

This part describes the client programming interfaces distributed with Postgres Pro. Each of these chapters can be read independently. There are many external programming interfaces for client programs that are distributed separately. They contain their own documentation ([Appendix I](#) lists some of the more popular ones). Readers of this part should be familiar with using SQL to manipulate and query the database (see [Part II](#)) and of course with the programming language of their choice.

Chapter 35. libpq — C Library

libpq is the C application programmer's interface to Postgres Pro. libpq is a set of library functions that allow client programs to pass queries to the Postgres Pro backend server and to receive the results of these queries.

libpq is also the underlying engine for several other Postgres Pro application interfaces, including those written for C++, Perl, Python, Tcl and ECPG. So some aspects of libpq's behavior will be important to you if you use one of those packages. In particular, [Section 35.15](#), [Section 35.16](#) and [Section 35.19](#) describe behavior that is visible to the user of any application that uses libpq.

Some short programs are included at the end of this chapter ([Section 35.22](#)) to show how to write programs that use libpq.

Client programs that use libpq must include the header file `libpq-fe.h` and must link with the libpq library.

35.1. Database Connection Control Functions

The following functions deal with making a connection to a Postgres Pro backend server. An application program can have several backend connections open at one time. (One reason to do that is to access more than one database.) Each connection is represented by a `PGconn` object, which is obtained from the function `PQconnectdb`, `PQconnectdbParams`, or `PQsetdbLogin`. Note that these functions will always return a non-null object pointer, unless perhaps there is too little memory even to allocate the `PGconn` object. The `PQstatus` function should be called to check the return value for a successful connection before queries are sent via the connection object.

Warning

If untrusted users have access to a database that has not adopted a [secure schema usage pattern](#), begin each session by removing publicly-writable schemas from `search_path`. One can set parameter key word `options` to value `-csearch_path=`. Alternately, one can issue `PQexec(conn, "SELECT pg_catalog.set_config('search_path', '', false)")` after connecting. This consideration is not specific to libpq; it applies to every interface for executing arbitrary SQL commands.

Warning

On Unix, forking a process with open libpq connections can lead to unpredictable results because the parent and child processes share the same sockets and operating system resources. For this reason, such usage is not recommended, though doing an `exec` from the child process to load a new executable is safe.

`PQconnectdbParams`

Makes a new connection to the database server.

```
PGconn *PQconnectdbParams(const char * const *keywords,
                          const char * const *values,
                          int expand_dbname);
```

This function opens a new database connection using the parameters taken from two `NULL`-terminated arrays. The first, `keywords`, is defined as an array of strings, each one being a key word. The second, `values`, gives the value for each key word. Unlike `PQsetdbLogin` below, the parameter set can be extended without changing the function signature, so use of this function (or its nonblocking analogs `PQconnectStartParams` and `PQconnectPoll`) is preferred for new application programming.

The currently recognized parameter key words are listed in [Section 35.1.2](#).

The passed arrays can be empty to use all default parameters, or can contain one or more parameter settings. They must be matched in length. Processing will stop at the first `NULL` entry in the `keywords` array. Also, if the `values` entry associated with a non-`NULL` `keywords` entry is `NULL` or an empty string, that entry is ignored and processing continues with the next pair of array entries.

When `expand_dbname` is non-zero, the value for the first `dbname` key word is checked to see if it is a *connection string*. If so, it is “expanded” into the individual connection parameters extracted from the string. The value is considered to be a connection string, rather than just a database name, if it contains an equal sign (=) or it begins with a URI scheme designator. (More details on connection string formats appear in [Section 35.1.1](#).) Only the first occurrence of `dbname` is treated in this way; any subsequent `dbname` parameter is processed as a plain database name.

In general the parameter arrays are processed from start to end. If any key word is repeated, the last value (that is not `NULL` or empty) is used. This rule applies in particular when a key word found in a connection string conflicts with one appearing in the `keywords` array. Thus, the programmer may determine whether array entries can override or be overridden by values taken from a connection string. Array entries appearing before an expanded `dbname` entry can be overridden by fields of the connection string, and in turn those fields are overridden by array entries appearing after `dbname` (but, again, only if those entries supply non-empty values).

After processing all the array entries and any expanded connection string, any connection parameters that remain unset are filled with default values. If an unset parameter's corresponding environment variable (see [Section 35.15](#)) is set, its value is used. If the environment variable is not set either, then the parameter's built-in default value is used.

`PQconnectdb`

Makes a new connection to the database server.

```
PGconn *PQconnectdb(const char *conninfo);
```

This function opens a new database connection using the parameters taken from the string `conninfo`.

The passed string can be empty to use all default parameters, or it can contain one or more parameter settings separated by whitespace, or it can contain a URI. See [Section 35.1.1](#) for details.

`PQsetdbLogin`

Makes a new connection to the database server.

```
PGconn *PQsetdbLogin(const char *pghost,
                    const char *pgport,
                    const char *pgoptions,
                    const char *pgtty,
                    const char *dbName,
                    const char *login,
                    const char *pwd);
```

This is the predecessor of `PQconnectdb` with a fixed set of parameters. It has the same functionality except that the missing parameters will always take on default values. Write `NULL` or an empty string for any one of the fixed parameters that is to be defaulted.

If the `dbName` contains an = sign or has a valid connection URI prefix, it is taken as a `conninfo` string in exactly the same way as if it had been passed to `PQconnectdb`, and the remaining parameters are then applied as specified for `PQconnectdbParams`.

`pgtty` is no longer used and any value passed will be ignored.

`PQsetdb`

Makes a new connection to the database server.

```
PGconn *PQsetdb(char *pghost,
```

```
char *pgport,  
char *pgoptions,  
char *pgtty,  
char *dbName);
```

This is a macro that calls [PQsetdbLogin](#) with null pointers for the *login* and *pwd* parameters. It is provided for backward compatibility with very old programs.

```
PQconnectStartParams  
PQconnectStart  
PQconnectPoll
```

Make a connection to the database server in a nonblocking manner.

```
PGconn *PQconnectStartParams(const char * const *keywords,  
                             const char * const *values,  
                             int expand_dbname);
```

```
PGconn *PQconnectStart(const char *conninfo);
```

```
PostgresPollingStatusType PQconnectPoll(PGconn *conn);
```

These three functions are used to open a connection to a database server such that your application's thread of execution is not blocked on remote I/O whilst doing so. The point of this approach is that the waits for I/O to complete can occur in the application's main loop, rather than down inside [PQconnectdbParams](#) or [PQconnectdb](#), and so the application can manage this operation in parallel with other activities.

With [PQconnectStartParams](#), the database connection is made using the parameters taken from the keywords and values arrays, and controlled by `expand_dbname`, as described above for [PQconnectdbParams](#).

With [PQconnectStart](#), the database connection is made using the parameters taken from the string `conninfo` as described above for [PQconnectdb](#).

Neither [PQconnectStartParams](#) nor [PQconnectStart](#) nor [PQconnectPoll](#) will block, so long as a number of restrictions are met:

- The `hostaddr` parameter must be used appropriately to prevent DNS queries from being made. See the documentation of this parameter in [Section 35.1.2](#) for details.
- If you call [PQtrace](#), ensure that the stream object into which you trace will not block.
- You must ensure that the socket is in the appropriate state before calling [PQconnectPoll](#), as described below.

To begin a nonblocking connection request, call [PQconnectStart](#) or [PQconnectStartParams](#). If the result is null, then libpq has been unable to allocate a new `PGconn` structure. Otherwise, a valid `PGconn` pointer is returned (though not yet representing a valid connection to the database). Next call `PQstatus(conn)`. If the result is `CONNECTION_BAD`, the connection attempt has already failed, typically because of invalid connection parameters.

If [PQconnectStart](#) or [PQconnectStartParams](#) succeeds, the next stage is to poll libpq so that it can proceed with the connection sequence. Use `PQsocket(conn)` to obtain the descriptor of the socket underlying the database connection. (Caution: do not assume that the socket remains the same across [PQconnectPoll](#) calls.) Loop thus: If [PQconnectPoll\(conn\)](#) last returned `PGRES_POLLING_READING`, wait until the socket is ready to read (as indicated by `select()`, `poll()`, or similar system function). Note that [PQsocketPoll](#) can help reduce boilerplate by abstracting the setup of `select(2)` or `poll(2)` if it is available on your system. Then call [PQconnectPoll\(conn\)](#) again. Conversely, if [PQconnectPoll\(conn\)](#) last returned `PGRES_POLLING_WRITING`, wait until the socket is ready to write, then call [PQconnectPoll\(conn\)](#) again. On the first iteration, i.e., if you have yet to call [PQconnectPoll](#), be-

have as if it last returned `PGRES_POLLING_WRITING`. Continue this loop until `PQconnectPoll(conn)` returns `PGRES_POLLING_FAILED`, indicating the connection procedure has failed, or `PGRES_POLLING_OK`, indicating the connection has been successfully made.

At any time during connection, the status of the connection can be checked by calling `PQstatus`. If this call returns `CONNECTION_BAD`, then the connection procedure has failed; if the call returns `CONNECTION_OK`, then the connection is ready. Both of these states are equally detectable from the return value of `PQconnectPoll`, described above. Other states might also occur during (and only during) an asynchronous connection procedure. These indicate the current stage of the connection procedure and might be useful to provide feedback to the user for example. These statuses are:

`CONNECTION_STARTED`

Waiting for connection to be made.

`CONNECTION_MADE`

Connection OK; waiting to send.

`CONNECTION_AWAITING_RESPONSE`

Waiting for a response from the server.

`CONNECTION_AUTH_OK`

Received authentication; waiting for backend start-up to finish.

`CONNECTION_SSL_STARTUP`

Negotiating SSL encryption.

`CONNECTION_GSS_STARTUP`

Negotiating GSS encryption.

`CONNECTION_CHECK_WRITABLE`

Checking if connection is able to handle write transactions.

`CONNECTION_CHECK_STANDBY`

Checking if connection is to a server in standby mode.

`CONNECTION_CONSUME`

Consuming any remaining response messages on connection.

Note that, although these constants will remain (in order to maintain compatibility), an application should never rely upon these occurring in a particular order, or at all, or on the status always being one of these documented values. An application might do something like this:

```
switch (PQstatus(conn))
{
    case CONNECTION_STARTED:
        feedback = "Connecting...";
        break;

    case CONNECTION_MADE:
        feedback = "Connected to server...";
        break;

    .
    .
    .

    default:
        feedback = "Connecting...";
}
```

The `connect_timeout` connection parameter is ignored when using `PQconnectPoll`; it is the application's responsibility to decide whether an excessive amount of time has elapsed. Otherwise, `PQconnectStart` followed by a `PQconnectPoll` loop is equivalent to `PQconnectdb`.

Note that when `PQconnectStart` or `PQconnectStartParams` returns a non-null pointer, you must call `PQfinish` when you are finished with it, in order to dispose of the structure and any associated memory blocks. This must be done even if the connection attempt fails or is abandoned.

`PQsocketPoll`

Poll a connection's underlying socket descriptor retrieved with `PQsocket`. The primary use of this function is iterating through the connection sequence described in the documentation of `PQconnectStartParams`.

```
typedef pg_int64 pg_usec_time_t;

int PQsocketPoll(int sock, int forRead, int forWrite,
                 pg_usec_time_t end_time);
```

This function performs polling of a file descriptor, optionally with a timeout. If `forRead` is nonzero, the function will terminate when the socket is ready for reading. If `forWrite` is nonzero, the function will terminate when the socket is ready for writing.

The timeout is specified by `end_time`, which is the time to stop waiting expressed as a number of microseconds since the Unix epoch (that is, `time_t` times 1 million). Timeout is infinite if `end_time` is `-1`. Timeout is immediate (no blocking) if `end_time` is `0` (or indeed, any time before now). Timeout values can be calculated conveniently by adding the desired number of microseconds to the result of `PQgetCurrentTimeUsec`. Note that the underlying system calls may have less than microsecond precision, so that the actual delay may be imprecise.

The function returns a value greater than `0` if the specified condition is met, `0` if a timeout occurred, or `-1` if an error occurred. The error can be retrieved by checking the `errno(3)` value. In the event both `forRead` and `forWrite` are zero, the function immediately returns a timeout indication.

`PQsocketPoll` is implemented using either `poll(2)` or `select(2)`, depending on platform. See `POLLIN` and `POLLOUT` from `poll(2)`, or `readfds` and `writefds` from `select(2)`, for more information.

`PQconnndefaults`

Returns the default connection options.

```
PQconninfoOption *PQconnndefaults(void);

typedef struct
{
    char    *keyword; /* The keyword of the option */
    char    *envvar;  /* Fallback environment variable name */
    char    *compiled; /* Fallback compiled in default value */
    char    *val;      /* Option's current value, or NULL */
    char    *label;    /* Label for field in connect dialog */
    char    *dispchar; /* Indicates how to display this field
                        in a connect dialog. Values are:
                        ""          Display entered value as is
                        "*"        Password field - hide value
                        "D"        Debug option - don't show by default */
    int     dispsize;  /* Field size in characters for dialog */
} PQconninfoOption;
```

Returns a connection options array. This can be used to determine all possible `PQconnectdb` options and their current default values. The return value points to an array of `PQconninfoOption` structures, which ends with an entry having a null keyword pointer. The null pointer is returned if memory

could not be allocated. Note that the current default values (`val` fields) will depend on environment variables and other context. A missing or invalid service file will be silently ignored. Callers must treat the connection options data as read-only.

After processing the options array, free it by passing it to [PQconninfoFree](#). If this is not done, a small amount of memory is leaked for each call to [PQconndefaults](#).

`PQconninfo`

Returns the connection options used by a live connection.

```
PQconninfoOption *PQconninfo(PGconn *conn);
```

Returns a connection options array. This can be used to determine all possible [PQconnectdb](#) options and the values that were used to connect to the server. The return value points to an array of `PQconninfoOption` structures, which ends with an entry having a null keyword pointer. All notes above for [PQconndefaults](#) also apply to the result of [PQconninfo](#).

`PQconninfoParse`

Returns parsed connection options from the provided connection string.

```
PQconninfoOption *PQconninfoParse(const char *conninfo, char **errmsg);
```

Parses a connection string and returns the resulting options as an array; or returns `NULL` if there is a problem with the connection string. This function can be used to extract the [PQconnectdb](#) options in the provided connection string. The return value points to an array of `PQconninfoOption` structures, which ends with an entry having a null keyword pointer.

All legal options will be present in the result array, but the `PQconninfoOption` for any option not present in the connection string will have `val` set to `NULL`; default values are not inserted.

If `errmsg` is not `NULL`, then `*errmsg` is set to `NULL` on success, else to a malloc'd error string explaining the problem. (It is also possible for `*errmsg` to be set to `NULL` and the function to return `NULL`; this indicates an out-of-memory condition.)

After processing the options array, free it by passing it to [PQconninfoFree](#). If this is not done, some memory is leaked for each call to [PQconninfoParse](#). Conversely, if an error occurs and `errmsg` is not `NULL`, be sure to free the error string using [PQfreemem](#).

`PQfinish`

Closes the connection to the server. Also frees memory used by the `PGconn` object.

```
void PQfinish(PGconn *conn);
```

Note that even if the server connection attempt fails (as indicated by [PQstatus](#)), the application should call [PQfinish](#) to free the memory used by the `PGconn` object. The `PGconn` pointer must not be used again after [PQfinish](#) has been called.

`PQreset`

Resets the communication channel to the server.

```
void PQreset(PGconn *conn);
```

This function will close the connection to the server and attempt to establish a new connection, using all the same parameters previously used. This might be useful for error recovery if a working connection is lost.

`PQresetStart`

`PQresetPoll`

Reset the communication channel to the server, in a nonblocking manner.

```
int PQresetStart(PGconn *conn);
```

```
PostgresPollingStatusType PQresetPoll(PGconn *conn);
```

These functions will close the connection to the server and attempt to establish a new connection, using all the same parameters previously used. This can be useful for error recovery if a working connection is lost. They differ from [PQreset](#) (above) in that they act in a nonblocking manner. These functions suffer from the same restrictions as [PQconnectStartParams](#), [PQconnectStart](#) and [PQconnectPoll](#).

To initiate a connection reset, call [PQresetStart](#). If it returns 0, the reset has failed. If it returns 1, poll the reset using [PQresetPoll](#) in exactly the same way as you would create the connection using [PQconnectPoll](#).

[PQpingParams](#)

[PQpingParams](#) reports the status of the server. It accepts connection parameters identical to those of [PQconnectdbParams](#), described above. It is not necessary to supply correct user name, password, or database name values to obtain the server status; however, if incorrect values are provided, the server will log a failed connection attempt.

```
PGPing PQpingParams(const char * const *keywords,
                    const char * const *values,
                    int expand_dbname);
```

The function returns one of the following values:

[PQPING_OK](#)

The server is running and appears to be accepting connections.

[PQPING_REJECT](#)

The server is running but is in a state that disallows connections (startup, shutdown, or crash recovery).

[PQPING_NO_RESPONSE](#)

The server could not be contacted. This might indicate that the server is not running, or that there is something wrong with the given connection parameters (for example, wrong port number), or that there is a network connectivity problem (for example, a firewall blocking the connection request).

[PQPING_NO_ATTEMPT](#)

No attempt was made to contact the server, because the supplied parameters were obviously incorrect or there was some client-side problem (for example, out of memory).

[PQping](#)

[PQping](#) reports the status of the server. It accepts connection parameters identical to those of [PQconnectdb](#), described above. It is not necessary to supply correct user name, password, or database name values to obtain the server status; however, if incorrect values are provided, the server will log a failed connection attempt.

```
PGPing PQping(const char *conninfo);
```

The return values are the same as for [PQpingParams](#).

[PQsetSSLKeyPassHook_OpenSSL](#)

[PQsetSSLKeyPassHook_OpenSSL](#) lets an application override libpq's [default handling of encrypted client certificate key files](#) using [sslpassword](#) or interactive prompting.

```
void PQsetSSLKeyPassHook_OpenSSL(PQsslKeyPassHook_OpenSSL_type hook);
```

The application passes a pointer to a callback function with signature:

```
int callback_fn(char *buf, int size, PGconn *conn);
```

which libpq will then call *instead of* its default `PQdefaultSSLKeyPassHook_OpenSSL` handler. The callback should determine the password for the key and copy it to result-buffer `buf` of size `size`. The string in `buf` must be null-terminated. The callback must return the length of the password stored in `buf` excluding the null terminator. On failure, the callback should set `buf[0] = '\0'` and return 0.

If the user specified an explicit key location, its path will be in `conn->sslkey` when the callback is invoked. This will be empty if the default key path is being used. For keys that are engine specifiers, it is up to engine implementations whether they use the OpenSSL password callback or define their own handling.

The app callback may choose to delegate unhandled cases to `PQdefaultSSLKeyPassHook_OpenSSL`, or call it first and try something else if it returns 0, or completely override it.

The callback *must not* escape normal flow control with exceptions, `longjmp(...)`, etc. It must return normally.

`PQgetSSLKeyPassHook_OpenSSL`

`PQgetSSLKeyPassHook_OpenSSL` returns the current client certificate key password hook, or NULL if none has been set.

```
PQsslKeyPassHook_OpenSSL_type PQgetSSLKeyPassHook_OpenSSL(void);
```

35.1.1. Connection Strings

Several libpq functions parse a user-specified string to obtain connection parameters. There are two accepted formats for these strings: plain keyword/value strings and URIs. URIs generally follow [RFC 3986](#), except that multi-host connection strings are allowed as further described below.

35.1.1.1. Keyword/Value Connection Strings

In the keyword/value format, each parameter setting is in the form `keyword = value`, with space(s) between settings. Spaces around a setting's equal sign are optional. To write an empty value, or a value containing spaces, surround it with single quotes, for example `keyword = 'a value'`. Single quotes and backslashes within a value must be escaped with a backslash, i.e., `\'` and `\\`.

Example:

```
host=localhost port=5432 dbname=mydb connect_timeout=10
```

The recognized parameter key words are listed in [Section 35.1.2](#).

35.1.1.2. Connection URIs

The general form for a connection URI is:

```
postgresql://[userspec@][hostspec][/dbname][?paramspec]
```

where `userspec` is:

```
user[:password]
```

and `hostspec` is:

```
[host][:port][,...]
```

and `paramspec` is:

`name=value[&...]`

The URI scheme designator can be either `postgresql://` or `postgres://`. Each of the remaining URI parts is optional. The following examples illustrate valid URI syntax:

```
postgresql://
postgresql://localhost
postgresql://localhost:5433
postgresql://localhost/mydb
postgresql://user@localhost
postgresql://user:secret@localhost
postgresql://other@localhost/otherdb?connect_timeout=10&application_name=myapp
postgresql://host1:123,host2:456/somedb?
hostorder=random&target_session_attrs=any&application_name=myapp
```

Values that would normally appear in the hierarchical part of the URI can alternatively be given as named parameters. For example:

```
postgresql:///mydb?host=localhost&port=5433
```

All named parameters must match key words listed in [Section 35.1.2](#), except that for compatibility with JDBC connection URIs, instances of `ssl=true` are translated into `sslmode=require`.

The connection URI needs to be encoded with [percent-encoding](#) if it includes symbols with special meaning in any of its parts. Here is an example where the equal sign (=) is replaced with `%3D` and the space character with `%20`:

```
postgresql://user@localhost:5433/mydb?options=-c%20synchronous_commit%3Doff
```

The host part may be either a host name or an IP address. To specify an IPv6 address, enclose it in square brackets:

```
postgresql://[2001:db8::1234]/database
```

The host part is interpreted as described for the parameter [host](#). In particular, a Unix-domain socket connection is chosen if the host part is either empty or looks like an absolute path name, otherwise a TCP/IP connection is initiated. Note, however, that the slash is a reserved character in the hierarchical part of the URI. So, to specify a non-standard Unix-domain socket directory, either omit the host part of the URI and specify the host as a named parameter, or percent-encode the path in the host part of the URI:

```
postgresql:///dbname?host=/var/lib/postgresql
postgresql://%2Fvar%2Flib%2Fpostgresql/dbname
```

It is possible to specify multiple host components, each with an optional port component, in a single URI. A URI of the form `postgresql://host1:port1,host2:port2,host3:port3/` is equivalent to a connection string of the form `host=host1,host2,host3 port=port1,port2,port3`. Each host will be tried in the order defined by the [hostorder](#) parameter until a connection is successfully established, all hosts have been tried, or [failover_timeout](#) is reached. For details, see [Section 35.1.1.3](#).

35.1.1.3. Specifying Multiple Hosts

It is possible to specify multiple hosts to connect to. In the Keyword/Value format, the `host`, `hostaddr`, and `port` options accept a comma-separated list of values. The same number of elements must be given in each option that is specified, such that e.g. the first `hostaddr` corresponds to the first host name, the second `hostaddr` corresponds to the second host name, and so forth. As an exception, if only one `port` is specified, it applies to all the hosts.

In the connection URI format, you can list multiple `host:port` pairs separated by commas in the `host` component of the URI.

In either format, a single host name can translate to multiple network addresses. A common example of this is a host that has both an IPv4 and an IPv6 address.

When multiple hosts are specified, or when a single host name is translated to multiple addresses, all the hosts and addresses will be tried in the order defined by the [hostorder](#) parameter until one succeeds. By default, each host is tried only once. You can allow multiple connection attempts to each host using the [failover_timeout](#) parameter, which specifies the maximum time for making a connection. If none of the hosts can be reached, the connection fails. If a connection is established successfully, but authentication fails, the remaining hosts in the list are not tried.

If a password file is used, you can have different passwords for different hosts. All the other connection options are the same for every host in the list; it is not possible to e.g., specify different usernames for different hosts.

35.1.2. Parameter Key Words

The currently recognized parameter key words are:

`host`

Name of host to connect to. If a host name looks like an absolute path name, it specifies Unix-domain communication rather than TCP/IP communication; the value is the name of the directory in which the socket file is stored. (On Unix, an absolute path name begins with a slash. On Windows, paths starting with drive letters are also recognized.) If the host name starts with @, it is taken as a Unix-domain socket in the abstract namespace (currently supported on Linux and Windows). The default behavior when `host` is not specified, or is empty, is to connect to a Unix-domain socket in `/tmp` (or whatever socket directory was specified when Postgres Pro was built). On Windows, the default is to connect to `localhost`.

A comma-separated list of host names is also accepted, in which case each host name in the list is tried in the order defined by the [hostorder](#) parameter; an empty item in the list selects the default behavior as explained above. See [Section 35.1.1.3](#) for details.

`hostaddr`

Numeric IP address of host to connect to. This should be in the standard IPv4 address format, e.g., `172.28.40.9`. If your machine supports IPv6, you can also use those addresses. TCP/IP communication is always used when a nonempty string is specified for this parameter. If this parameter is not specified, the value of `host` will be looked up to find the corresponding IP address — or, if `host` specifies an IP address, that value will be used directly.

Using `hostaddr` allows the application to avoid a host name look-up, which might be important in applications with time constraints. However, a host name is required for GSSAPI or SSPI authentication methods, as well as for `verify-full` SSL certificate verification. The following rules are used:

- If `host` is specified without `hostaddr`, a host name lookup occurs. (When using `PQconnectPoll`, the lookup occurs when `PQconnectPoll` first considers this host name, and it may cause `PQconnectPoll` to block for a significant amount of time.)
- If `hostaddr` is specified without `host`, the value for `hostaddr` gives the server network address. The connection attempt will fail if the authentication method requires a host name.
- If both `host` and `hostaddr` are specified, the value for `hostaddr` gives the server network address. The value for `host` is ignored unless the authentication method requires it, in which case it will be used as the host name.

Note that authentication is likely to fail if `host` is not the name of the server at network address `hostaddr`. Also, when both `host` and `hostaddr` are specified, `host` is used to identify the connection in a password file (see [Section 35.16](#)).

A comma-separated list of `hostaddr` values is also accepted, in which case each host in the list is tried in the order defined by the [hostorder](#) parameter. An empty item in the list causes the corresponding host name to be used, or the default host name if that is empty as well. See [Section 35.1.1.3](#) for details.

Without either a host name or host address, libpq will connect using a local Unix-domain socket; or on Windows, it will attempt to connect to `localhost`.

hostorder

Defines the order of connecting to the hosts specified by the [host](#) or [hostaddr](#) parameter.

If the value of this parameter is `sequential` (the default), connections to the hosts are attempted in the order in which the hosts are listed.

If the value is `random`, the host to connect to is randomly picked from the list. This setting can help you balance the load between several cluster nodes instead of always connecting to the first available host in the list. You can also use the [target_session_attrs](#) parameter to forbid read-only connections.

port

Port number to connect to at the server host, or socket file name extension for Unix-domain connections. If multiple hosts were given in the [host](#) or [hostaddr](#) parameters, this parameter may specify a comma-separated list of ports of the same length as the host list, or it may specify a single port number to be used for all hosts. An empty string, or an empty item in a comma-separated list, specifies the default port number established when Postgres Pro was built.

dbname

The database name. Defaults to be the same as the user name. In certain contexts, the value is checked for extended formats; see [Section 35.1.1](#) for more details on those.

user

Postgres Pro user name to connect as. Defaults to be the same as the operating system name of the user running the application.

password

Password to be used if the server demands password authentication.

passfile

Specifies the name of the file used to store passwords (see [Section 35.16](#)). Defaults to `~/.pgpass`, or `%APPDATA%\postgresql\pgpass.conf` on Microsoft Windows. (No error is reported if this file does not exist.)

require_auth

Specifies the authentication method that the client requires from the server. If the server does not use the required method to authenticate the client, or if the authentication handshake is not fully completed by the server, the connection will fail. A comma-separated list of methods may also be provided, of which the server must use exactly one in order for the connection to succeed. By default, any authentication method is accepted, and the server is free to skip authentication altogether.

Methods may be negated with the addition of a `!` prefix, in which case the server must *not* attempt the listed method; any other method is accepted, and the server is free not to authenticate the client at all. If a comma-separated list is provided, the server may not attempt *any* of the listed negated methods. Negated and non-negated forms may not be combined in the same setting.

As a final special case, the `none` method requires the server not to use an authentication challenge. (It may also be negated, to require some form of authentication.)

The following methods may be specified:

password

The server must request plaintext password authentication.

md5

The server must request MD5 hashed password authentication.

gss

The server must either request a Kerberos handshake via GSSAPI or establish a GSS-encrypted channel (see also [gssencmode](#)).

sspi

The server must request Windows SSPI authentication.

scram-sha-256

The server must successfully complete a SCRAM-SHA-256 authentication exchange with the client.

none

The server must not prompt the client for an authentication exchange. (This does not prohibit client certificate authentication via TLS, nor GSS authentication via its encrypted transport.)

channel_binding

This option controls the client's use of channel binding. A setting of `require` means that the connection must employ channel binding, `prefer` means that the client will choose channel binding if available, and `disable` prevents the use of channel binding. The default is `prefer` if Postgres Pro is compiled with SSL support; otherwise the default is `disable`.

Channel binding is a method for the server to authenticate itself to the client. It is only supported over SSL connections with Postgres Pro 11 or later servers using the `SCRAM` authentication method.

reusepass

Specifies whether reusing the prompted password is allowed when reestablishing connections. To forbid password reuse, set `reusepass` to 0. You can use this option together with [idle_session_timeout](#) to define the reconnection behavior for sessions terminated by timeout.

Default: 1

connect_timeout

Maximum time to wait while connecting, in seconds (write as a decimal integer, e.g., 10). Zero, negative, or not specified means wait indefinitely. This timeout applies separately to each host name or IP address. For example, if you specify two hosts and `connect_timeout` is 5, each host will time out if no connection is made within 5 seconds, so the total time spent waiting for a connection might be up to 10 seconds.

failover_timeout

Maximum time, in seconds, to cyclically retry to connect to the hosts specified in the connection string. If this parameter is not set (the default), connection to each host is tried only once. The specified value must be a decimal integer number.

When using this parameter in clusters with primary-standby replication, make sure that its value provides enough time for a standby node to become the new primary if the current primary node fails. In this case, clients will be able to reconnect to the new primary once failover is complete.

client_encoding

This sets the `client_encoding` configuration parameter for this connection. In addition to the values accepted by the corresponding server option, you can use `auto` to determine the right encoding from the current locale in the client (`LC_CTYPE` environment variable on Unix systems).

options

Specifies command-line options to send to the server at connection start. For example, setting this to `-c geqo=off` or `--geqo=off` sets the session's value of the `geqo` parameter to `off`. Spaces within this string are considered to separate command-line arguments, unless escaped with a backslash (`\`); write `\\` to represent a literal backslash. For a detailed discussion of the available options, consult [Chapter 19](#).

application_name

Specifies a value for the [application_name](#) configuration parameter.

fallback_application_name

Specifies a fallback value for the [application_name](#) configuration parameter. This value will be used if no value has been given for `application_name` via a connection parameter or the `PGAPPNAME` environment variable. Specifying a fallback name is useful in generic utility programs that wish to set a default application name but allow it to be overridden by the user.

keepalives

Controls whether client-side TCP keepalives are used. The default value is 1, meaning on, but you can change this to 0, meaning off, if keepalives are not wanted. This parameter is ignored for connections made via a Unix-domain socket.

keepalives_idle

Controls the number of seconds of inactivity after which TCP should send a keepalive message to the server. A value of zero uses the system default. This parameter is ignored for connections made via a Unix-domain socket, or if keepalives are disabled. It is only supported on systems where `TCP_KEEPIDLE` or an equivalent socket option is available, and on Windows; on other systems, it has no effect.

keepalives_interval

Controls the number of seconds after which a TCP keepalive message that is not acknowledged by the server should be retransmitted. A value of zero uses the system default. This parameter is ignored for connections made via a Unix-domain socket, or if keepalives are disabled. It is only supported on systems where `TCP_KEEPINTVL` or an equivalent socket option is available, and on Windows; on other systems, it has no effect.

keepalives_count

Controls the number of TCP keepalives that can be lost before the client's connection to the server is considered dead. A value of zero uses the system default. This parameter is ignored for connections made via a Unix-domain socket, or if keepalives are disabled. It is only supported on systems where `TCP_KEEPCNT` or an equivalent socket option is available; on other systems, it has no effect.

tcp_user_timeout

Controls the number of milliseconds that transmitted data may remain unacknowledged before a connection is forcibly closed. A value of zero uses the system default. This parameter is ignored for connections made via a Unix-domain socket. It is only supported on systems where `TCP_USER_TIMEOUT` is available; on other systems, it has no effect.

replication

This option determines whether the connection should use the replication protocol instead of the normal protocol. This is what Postgres Pro replication connections as well as tools such as `pg_basebackup` use internally, but it can also be used by third-party applications. For a description of the replication protocol, consult [Section 56.4](#).

The following values, which are case-insensitive, are supported:

true, on, yes, 1

The connection goes into physical replication mode.

database

The connection goes into logical replication mode, connecting to the database specified in the `dbname` parameter.

false, off, no, 0

The connection is a regular one, which is the default behavior.

In physical or logical replication mode, only the simple query protocol can be used.

gssencmode

This option determines whether or with what priority a secure GSS TCP/IP connection will be negotiated with the server. There are three modes:

disable

only try a non-GSSAPI-encrypted connection

prefer (default)

if there are GSSAPI credentials present (i.e., in a credentials cache), first try a GSSAPI-encrypted connection; if that fails or there are no credentials, try a non-GSSAPI-encrypted connection. This is the default when Postgres Pro has been compiled with GSSAPI support.

require

only try a GSSAPI-encrypted connection

`gssencmode` is ignored for Unix domain socket communication. If Postgres Pro is compiled without GSSAPI support, using the `require` option will cause an error, while `prefer` will be accepted but libpq will not actually attempt a GSSAPI-encrypted connection.

sslmode

This option determines whether or with what priority a secure SSL TCP/IP connection will be negotiated with the server. There are six modes:

disable

only try a non-SSL connection

allow

first try a non-SSL connection; if that fails, try an SSL connection

prefer (default)

first try an SSL connection; if that fails, try a non-SSL connection

require

only try an SSL connection. If a root CA file is present, verify the certificate in the same way as if `verify-ca` was specified

verify-ca

only try an SSL connection, and verify that the server certificate is issued by a trusted certificate authority (CA)

verify-full

only try an SSL connection, verify that the server certificate is issued by a trusted CA and that the requested server host name matches that in the certificate

See [Section 35.19](#) for a detailed description of how these options work.

`sslmode` is ignored for Unix domain socket communication. If Postgres Pro is compiled without SSL support, using options `require`, `verify-ca`, or `verify-full` will cause an error, while options `allow` and `prefer` will be accepted but libpq will not actually attempt an SSL connection.

Note that if GSSAPI encryption is possible, that will be used in preference to SSL encryption, regardless of the value of `sslmode`. To force use of SSL encryption in an environment that has working GSSAPI infrastructure (such as a Kerberos server), also set `gssencmode` to `disable`.

`requiressl`

This option is deprecated in favor of the `sslmode` setting.

If set to 1, an SSL connection to the server is required (this is equivalent to `sslmode require`). libpq will then refuse to connect if the server does not accept an SSL connection. If set to 0 (default), libpq will negotiate the connection type with the server (equivalent to `sslmode prefer`). This option is only available if Postgres Pro is compiled with SSL support.

`sslnegotiation`

This option controls how SSL encryption is negotiated with the server, if SSL is used. In the default `postgres` mode, the client first asks the server if SSL is supported. In `direct` mode, the client starts the standard SSL handshake directly after establishing the TCP/IP connection. Traditional Postgres Pro protocol negotiation is the most flexible with different server configurations. If the server is known to support direct SSL connections then the latter requires one fewer round trip reducing connection latency and also allows the use of protocol agnostic SSL network tools. The direct SSL option was introduced in Postgres Pro version 17.

`postgres`

perform Postgres Pro protocol negotiation. This is the default if the option is not provided.

`direct`

start SSL handshake directly after establishing the TCP/IP connection. This is only allowed with `sslmode=require` or higher, because the weaker settings could lead to unintended fallback to plaintext authentication when the server does not support direct SSL handshake.

`sslcompression`

If set to 1, data sent over SSL connections will be compressed. If set to 0, compression will be disabled. The default is 0. This parameter is ignored if a connection without SSL is made.

SSL compression is nowadays considered insecure and its use is no longer recommended. OpenSSL 1.1.0 disabled compression by default, and many operating system distributions disabled it in prior versions as well, so setting this parameter to on will not have any effect if the server does not accept compression. Postgres Pro 14 disabled compression completely in the backend.

If security is not a primary concern, compression can improve throughput if the network is the bottleneck. Disabling compression can improve response time and throughput if CPU performance is the limiting factor.

`sslcert`

This parameter specifies the file name of the client SSL certificate, replacing the default `~/.postgresql/postgresql.crt`. This parameter is ignored if an SSL connection is not made.

`sslkey`

This parameter specifies the location for the secret key used for the client certificate. It can either specify a file name that will be used instead of the default `~/.postgresql/postgresql.key`, or it

can specify a key obtained from an external “engine” (engines are OpenSSL loadable modules). An external engine specification should consist of a colon-separated engine name and an engine-specific key identifier. This parameter is ignored if an SSL connection is not made.

sslpassword

This parameter specifies the password for the secret key specified in `sslkey`, allowing client certificate private keys to be stored in encrypted form on disk even when interactive passphrase input is not practical.

Specifying this parameter with any non-empty value suppresses the `Enter PEM pass phrase: prompt` that OpenSSL will emit by default when an encrypted client certificate key is provided to libpq.

If the key is not encrypted this parameter is ignored. The parameter has no effect on keys specified by OpenSSL engines unless the engine uses the OpenSSL password callback mechanism for prompts.

There is no environment variable equivalent to this option, and no facility for looking it up in `.pgpass`. It can be used in a service file connection definition. Users with more sophisticated uses should consider using OpenSSL engines and tools like PKCS#11 or USB crypto offload devices.

sslcertmode

This option determines whether a client certificate may be sent to the server, and whether the server is required to request one. There are three modes:

disable

A client certificate is never sent, even if one is available (default location or provided via [sslcert](#)).

allow (default)

A certificate may be sent, if the server requests one and the client has one to send.

require

The server *must* request a certificate. The connection will fail if the client does not send a certificate and the server successfully authenticates the client anyway.

Note

`sslcertmode=require` doesn't add any additional security, since there is no guarantee that the server is validating the certificate correctly; Postgres Pro servers generally request TLS certificates from clients whether they validate them or not. The option may be useful when troubleshooting more complicated TLS setups.

sslrootcert

This parameter specifies the name of a file containing SSL certificate authority (CA) certificate(s). If the file exists, the server's certificate will be verified to be signed by one of these authorities. The default is `~/.postgresql/root.crt`.

The special value `system` may be specified instead, in which case the trusted CA roots from the SSL implementation will be loaded. The exact locations of these root certificates differ by SSL implementation and platform. For OpenSSL in particular, the locations may be further modified by the `SSL_CERT_DIR` and `SSL_CERT_FILE` environment variables.

Note

When using `sslrootcert=system`, the default `sslmode` is changed to `verify-full`, and any weaker setting will result in an error. In most cases it is trivial for anyone to obtain a certificate

trusted by the system for a hostname they control, rendering `verify-ca` and all weaker modes useless.

The magic `system` value will take precedence over a local certificate file with the same name. If for some reason you find yourself in this situation, use an alternative path like `ssl-rootcert=./system` instead.

`sslcr1`

This parameter specifies the file name of the SSL server certificate revocation list (CRL). Certificates listed in this file, if it exists, will be rejected while attempting to authenticate the server's certificate. If neither `sslcr1` nor `sslcrldir` is set, this setting is taken as `~/.postgresql/root.crl`.

`sslcrldir`

This parameter specifies the directory name of the SSL server certificate revocation list (CRL). Certificates listed in the files in this directory, if it exists, will be rejected while attempting to authenticate the server's certificate.

The directory needs to be prepared with the OpenSSL command `openssl rehash` or `c_rehash`. See its documentation for details.

Both `sslcr1` and `sslcrldir` can be specified together.

`sslsni`

If set to 1 (default), libpq sets the TLS extension “Server Name Indication” (SNI) on SSL-enabled connections. By setting this parameter to 0, this is turned off.

The Server Name Indication can be used by SSL-aware proxies to route connections without having to decrypt the SSL stream. (Note that unless the proxy is aware of the Postgres Pro protocol handshake this would require setting `sslnegotiation` to `direct`.) However, SNI makes the destination host name appear in cleartext in the network traffic, so it might be undesirable in some cases.

`requirepeer`

This parameter specifies the operating-system user name of the server, for example `requirepeer=postgres`. When making a Unix-domain socket connection, if this parameter is set, the client checks at the beginning of the connection that the server process is running under the specified user name; if it is not, the connection is aborted with an error. This parameter can be used to provide server authentication similar to that available with SSL certificates on TCP/IP connections. (Note that if the Unix-domain socket is in `/tmp` or another publicly writable location, any user could start a server listening there. Use this parameter to ensure that you are connected to a server run by a trusted user.) This option is only supported on platforms for which the `peer` authentication method is implemented; see [Section 20.9](#).

`ssl_min_protocol_version`

This parameter specifies the minimum SSL/TLS protocol version to allow for the connection. Valid values are `TLSv1`, `TLSv1.1`, `TLSv1.2` and `TLSv1.3`. The supported protocols depend on the version of OpenSSL used, older versions not supporting the most modern protocol versions. If not specified, the default is `TLSv1.2`, which satisfies industry best practices as of this writing.

`ssl_max_protocol_version`

This parameter specifies the maximum SSL/TLS protocol version to allow for the connection. Valid values are `TLSv1`, `TLSv1.1`, `TLSv1.2` and `TLSv1.3`. The supported protocols depend on the version of OpenSSL used, older versions not supporting the most modern protocol versions. If not set, this parameter is ignored and the connection will use the maximum bound defined by the backend, if set. Setting the maximum protocol version is mainly useful for testing or if some component has issues working with a newer protocol.

krbsrvname

Kerberos service name to use when authenticating with GSSAPI. This must match the service name specified in the server configuration for Kerberos authentication to succeed. (See also [Section 20.6](#).) The default value is normally `postgres`, but that can be changed when building Postgres Pro via the `--with-krb-srvnam` option of `configure`. In most environments, this parameter never needs to be changed. Some Kerberos implementations might require a different service name, such as Microsoft Active Directory which requires the service name to be in upper case (`POSTGRES`).

gsslib

GSS library to use for GSSAPI authentication. Currently this is disregarded except on Windows builds that include both GSSAPI and SSPI support. In that case, set this to `gssapi` to cause libpq to use the GSSAPI library for authentication instead of the default SSPI.

gssdelegation

Forward (delegate) GSS credentials to the server. The default is `0` which means credentials will not be forwarded to the server. Set this to `1` to have credentials forwarded when possible.

service

Service name to use for additional parameters. It specifies a service name in `pg_service.conf` that holds additional connection parameters. This allows applications to specify only a service name so connection parameters can be centrally maintained. See [Section 35.17](#).

target_session_attrs

This option determines whether the session must have certain properties to be acceptable. It's typically used in combination with multiple host names to select the first acceptable alternative among several hosts. There are six modes:

any (default)

any successful connection is acceptable

read-write

session must accept read-write transactions by default (that is, the server must not be in hot standby mode and the `default_transaction_read_only` parameter must be `off`)

read-only

session must not accept read-write transactions by default (the converse)

primary

server must not be in hot standby mode

standby

server must be in hot standby mode

prefer-standby

first try to find a standby server, but if none of the listed hosts is a standby server, try again in any mode

load_balance_hosts

Controls the order in which the client tries to connect to the available hosts and addresses. Once a connection attempt is successful no other hosts and addresses will be tried. This parameter is typically used in combination with multiple host names or a DNS record that returns multiple IPs. This parameter can be used in combination with [target_session_attrs](#) to, for example, load balance over

standby servers only. Once successfully connected, subsequent queries on the returned connection will all be sent to the same server. There are currently two modes:

`disable` (default)

No load balancing across hosts is performed. Hosts are tried in the order in which they are provided and addresses are tried in the order they are received from DNS or a hosts file.

`random`

Hosts and addresses are tried in random order. This value is mostly useful when opening multiple connections at the same time, possibly from different machines. This way connections can be load balanced across multiple Postgres Pro servers.

While random load balancing, due to its random nature, will almost never result in a completely uniform distribution, it statistically gets quite close. One important aspect here is that this algorithm uses two levels of random choices: First the hosts will be resolved in random order. Then secondly, before resolving the next host, all resolved addresses for the current host will be tried in random order. This behaviour can skew the amount of connections each node gets greatly in certain cases, for instance when some hosts resolve to more addresses than others. But such a skew can also be used on purpose, e.g. to increase the number of connections a larger server gets by providing its hostname multiple times in the host string.

When using this value it's recommended to also configure a reasonable value for [connect_timeout](#). Because then, if one of the nodes that are used for load balancing is not responding, a new node will be tried.

`target_server_type`

Same as [target_session_attrs](#). This parameter is provided for backward compatibility with Postgres Pro Shardman 9.6.

35.2. Connection Status Functions

These functions can be used to interrogate the status of an existing database connection object.

Tip

libpq application programmers should be careful to maintain the `PGconn` abstraction. Use the accessor functions described below to get at the contents of `PGconn`. Reference to internal `PGconn` fields using `libpq-int.h` is not recommended because they are subject to change in the future.

The following functions return parameter values established at connection. These values are fixed for the life of the connection. If a multi-host connection string is used, the values of [PQhost](#), [PQport](#), and [PQpass](#) can change if a new connection is established using the same `PGconn` object. Other values are fixed for the lifetime of the `PGconn` object.

`PQdb`

Returns the database name of the connection.

```
char *PQdb(const PGconn *conn);
```

`PQuser`

Returns the user name of the connection.

```
char *PQuser(const PGconn *conn);
```

`PQpass`

Returns the password of the connection.


```
char *PQpass(const PGconn *conn);
```

`PQpass` will return either the password specified in the connection parameters, or if there was none and the password was obtained from the [password file](#), it will return that. In the latter case, if multiple hosts were specified in the connection parameters, it is not possible to rely on the result of `PQpass` until the connection is established. The status of the connection can be checked using the function `PQstatus`.

PQhost

Returns the server host name of the active connection. This can be a host name, an IP address, or a directory path if the connection is via Unix socket. (The path case can be distinguished because it will always be an absolute path, beginning with `/`.)

```
char *PQhost(const PGconn *conn);
```

If the connection parameters specified both `host` and `hostaddr`, then `PQhost` will return the `host` information. If only `hostaddr` was specified, then that is returned. If multiple hosts were specified in the connection parameters, `PQhost` returns the host actually connected to.

`PQhost` returns `NULL` if the `conn` argument is `NULL`. Otherwise, if there is an error producing the host information (perhaps if the connection has not been fully established or there was an error), it returns an empty string.

If multiple hosts were specified in the connection parameters, it is not possible to rely on the result of `PQhost` until the connection is established. The status of the connection can be checked using the function `PQstatus`.

PQhostaddr

Returns the server IP address of the active connection. This can be the address that a host name resolved to, or an IP address provided through the `hostaddr` parameter.

```
char *PQhostaddr(const PGconn *conn);
```

`PQhostaddr` returns `NULL` if the `conn` argument is `NULL`. Otherwise, if there is an error producing the host information (perhaps if the connection has not been fully established or there was an error), it returns an empty string.

PQport

Returns the port of the active connection.

```
char *PQport(const PGconn *conn);
```

If multiple ports were specified in the connection parameters, `PQport` returns the port actually connected to.

`PQport` returns `NULL` if the `conn` argument is `NULL`. Otherwise, if there is an error producing the port information (perhaps if the connection has not been fully established or there was an error), it returns an empty string.

If multiple ports were specified in the connection parameters, it is not possible to rely on the result of `PQport` until the connection is established. The status of the connection can be checked using the function `PQstatus`.

PQtty

This function no longer does anything, but it remains for backwards compatibility. The function always return an empty string, or `NULL` if the `conn` argument is `NULL`.

```
char *PQtty(const PGconn *conn);
```

PQoptions

Returns the command-line options passed in the connection request.

```
char *PQoptions(const PGconn *conn);
```

The following functions return status data that can change as operations are executed on the `PGconn` object.

PQstatus

Returns the status of the connection.

```
ConnStatusType PQstatus(const PGconn *conn);
```

The status can be one of a number of values. However, only two of these are seen outside of an asynchronous connection procedure: `CONNECTION_OK` and `CONNECTION_BAD`. A good connection to the database has the status `CONNECTION_OK`. A failed connection attempt is signaled by status `CONNECTION_BAD`. Ordinarily, an OK status will remain so until [PQfinish](#), but a communications failure might result in the status changing to `CONNECTION_BAD` prematurely. In that case the application could try to recover by calling [PQreset](#).

See the entry for [PQconnectStartParams](#), `PQconnectStart` and `PQconnectPoll` with regards to other status codes that might be returned.

PQtransactionStatus

Returns the current in-transaction status of the server.

```
PGTransactionStatusType PQtransactionStatus(const PGconn *conn);
```

The status can be `PQTRANS_IDLE` (currently idle), `PQTRANS_ACTIVE` (a command is in progress), `PQTRANS_INTRANS` (idle, in a valid transaction block), or `PQTRANS_INERROR` (idle, in a failed transaction block). `PQTRANS_UNKNOWN` is reported if the connection is bad. `PQTRANS_ACTIVE` is reported only when a query has been sent to the server and not yet completed.

PQparameterStatus

Looks up a current parameter setting of the server.

```
const char *PQparameterStatus(const PGconn *conn, const char *paramName);
```

Certain parameter values are reported by the server automatically at connection startup or whenever their values change. [PQparameterStatus](#) can be used to interrogate these settings. It returns the current value of a parameter if known, or `NULL` if the parameter is not known.

Parameters reported as of the current release include:

<code>application_name</code>	<code>is_superuser</code>
<code>client_encoding</code>	<code>scram_iterations</code>
<code>DateStyle</code>	<code>server_encoding</code>
<code>default_transaction_read_only</code>	<code>server_version</code>
<code>in_hot_standby</code>	<code>session_authorization</code>
<code>integer_datetimes</code>	<code>standard_conforming_strings</code>
<code>IntervalStyle</code>	<code>TimeZone</code>

(`default_transaction_read_only` and `in_hot_standby` were not reported by releases before 14; `scram_iterations` was not reported by releases before 16.) Note that `server_version`, `server_encoding` and `integer_datetimes` cannot change after startup.

If no value for `standard_conforming_strings` is reported, applications can assume it is off, that is, backslashes are treated as escapes in string literals. Also, the presence of this parameter can be taken as an indication that the escape string syntax (`E' ... '`) is accepted.

Although the returned pointer is declared `const`, it in fact points to mutable storage associated with the `PGconn` structure. It is unwise to assume the pointer will remain valid across queries.

`PQprotocolVersion`

Interrogates the frontend/backend protocol being used.

```
int PQprotocolVersion(const PGconn *conn);
```

Applications might wish to use this function to determine whether certain features are supported. Currently, the possible values are 3 (3.0 protocol), or zero (connection bad). The protocol version will not change after connection startup is complete, but it could theoretically change during a connection reset. The 3.0 protocol is supported by PostgreSQL server versions 7.4 and above.

`PQserverVersion`

Returns an integer representing the server version.

```
int PQserverVersion(const PGconn *conn);
```

Applications might use this function to determine the version of the database server they are connected to. The result is formed by multiplying the server's major version number by 10000 and adding the minor version number. For example, version 10.1 will be returned as 100001, and version 11.0 will be returned as 110000. Zero is returned if the connection is bad.

Prior to major version 10, Postgres Pro used three-part version numbers in which the first two parts together represented the major version. For those versions, `PQserverVersion` uses two digits for each part; for example version 9.1.5 will be returned as 90105, and version 9.2.0 will be returned as 90200.

Therefore, for purposes of determining feature compatibility, applications should divide the result of `PQserverVersion` by 100 not 10000 to determine a logical major version number. In all release series, only the last two digits differ between minor releases (bug-fix releases).

`PQerrorMessage`

Returns the error message most recently generated by an operation on the connection.

```
char *PQerrorMessage(const PGconn *conn);
```

Nearly all libpq functions will set a message for `PQerrorMessage` if they fail. Note that by libpq convention, a nonempty `PQerrorMessage` result can consist of multiple lines, and will include a trailing newline. The caller should not free the result directly. It will be freed when the associated `PGconn` handle is passed to `PQfinish`. The result string should not be expected to remain the same across operations on the `PGconn` structure.

`PQsocket`

Obtains the file descriptor number of the connection socket to the server. A valid descriptor will be greater than or equal to 0; a result of -1 indicates that no server connection is currently open. (This will not change during normal operation, but could change during connection setup or reset.)

```
int PQsocket(const PGconn *conn);
```

`PQbackendPID`

Returns the process ID (PID) of the backend process handling this connection.

```
int PQbackendPID(const PGconn *conn);
```

The backend PID is useful for debugging purposes and for comparison to `NOTIFY` messages (which include the PID of the notifying backend process). Note that the PID belongs to a process executing on the database server host, not the local host!

PQconnectionNeedsPassword

Returns true (1) if the connection authentication method required a password, but none was available. Returns false (0) if not.

```
int PQconnectionNeedsPassword(const PGconn *conn);
```

This function can be applied after a failed connection attempt to decide whether to prompt the user for a password.

PQconnectionUsedPassword

Returns true (1) if the connection authentication method used a password. Returns false (0) if not.

```
int PQconnectionUsedPassword(const PGconn *conn);
```

This function can be applied after either a failed or successful connection attempt to detect whether the server demanded a password.

PQconnectionUsedGSSAPI

Returns true (1) if the connection authentication method used GSSAPI. Returns false (0) if not.

```
int PQconnectionUsedGSSAPI(const PGconn *conn);
```

This function can be applied to detect whether the connection was authenticated with GSSAPI.

The following functions return information related to SSL. This information usually doesn't change after a connection is established.

PQsslInUse

Returns true (1) if the connection uses SSL, false (0) if not.

```
int PQsslInUse(const PGconn *conn);
```

PQsslAttribute

Returns SSL-related information about the connection.

```
const char *PQsslAttribute(const PGconn *conn, const char *attribute_name);
```

The list of available attributes varies depending on the SSL library being used and the type of connection. Returns NULL if the connection does not use SSL or the specified attribute name is not defined for the library in use.

The following attributes are commonly available:

library

Name of the SSL implementation in use. (Currently, only "OpenSSL" is implemented)

protocol

SSL/TLS version in use. Common values are "TLSv1", "TLSv1.1" and "TLSv1.2", but an implementation may return other strings if some other protocol is used.

key_bits

Number of key bits used by the encryption algorithm.

cipher

A short name of the ciphersuite used, e.g., "DHE-RSA-DES-CBC3-SHA". The names are specific to each SSL implementation.

compression

Returns "on" if SSL compression is in use, else it returns "off".

alpn

Application protocol selected by the TLS Application-Layer Protocol Negotiation (ALPN) extension. The only protocol supported by libpq is `postgresql`, so this is mainly useful for checking whether the server supported ALPN or not. Empty string if ALPN was not used.

As a special case, the `library` attribute may be queried without a connection by passing `NULL` as the `conn` argument. The result will be the default SSL library name, or `NULL` if libpq was compiled without any SSL support. (Prior to Postgres Pro version 15, passing `NULL` as the `conn` argument always resulted in `NULL`. Client programs needing to differentiate between the newer and older implementations of this case may check the `LIBPQ_HAS_SSL_LIBRARY_DETECTION` feature macro.)

PQsslAttributeNames

Returns an array of SSL attribute names that can be used in `PQsslAttribute()`. The array is terminated by a `NULL` pointer.

```
const char * const * PQsslAttributeNames(const PGconn *conn);
```

If `conn` is `NULL`, the attributes available for the default SSL library are returned, or an empty list if libpq was compiled without any SSL support. If `conn` is not `NULL`, the attributes available for the SSL library in use for the connection are returned, or an empty list if the connection is not encrypted.

PQsslStruct

Returns a pointer to an SSL-implementation-specific object describing the connection. Returns `NULL` if the connection is not encrypted or the requested type of object is not available from the connection's SSL implementation.

```
void *PQsslStruct(const PGconn *conn, const char *struct_name);
```

The struct(s) available depend on the SSL implementation in use. For OpenSSL, there is one struct, available under the name `OpenSSL`, and it returns a pointer to OpenSSL's SSL struct. To use this function, code along the following lines could be used:

```
#include <libpq-fe.h>
#include <openssl/ssl.h>

...

SSL *ssl;

dbconn = PQconnectdb(...);
...

ssl = PQsslStruct(dbconn, "OpenSSL");
if (ssl)
{
    /* use OpenSSL functions to access ssl */
}
```

This structure can be used to verify encryption levels, check server certificates, and more. Refer to the OpenSSL documentation for information about this structure.

PQgetssl

Returns the SSL structure used in the connection, or `NULL` if SSL is not in use.

```
void *PQgetssl(const PGconn *conn);
```

This function is equivalent to `PQsslStruct(conn, "OpenSSL")`. It should not be used in new applications, because the returned struct is specific to OpenSSL and will not be available if another SSL

implementation is used. To check if a connection uses SSL, call [PQsslInUse](#) instead, and for more details about the connection, use [PQsslAttribute](#).

35.3. Command Execution Functions

Once a connection to a database server has been successfully established, the functions described here are used to perform SQL queries and commands.

35.3.1. Main Functions

[PQexec](#)

Submits a command to the server and waits for the result.

```
PGresult *PQexec(PGconn *conn, const char *command);
```

Returns a `PGresult` pointer or possibly a null pointer. A non-null pointer will generally be returned except in out-of-memory conditions or serious errors such as inability to send the command to the server. The [PQresultStatus](#) function should be called to check the return value for any errors (including the value of a null pointer, in which case it will return `PGRES_FATAL_ERROR`). Use [PQerrorMessage](#) to get more information about such errors.

The command string can include multiple SQL commands (separated by semicolons). Multiple queries sent in a single [PQexec](#) call are processed in a single transaction, unless there are explicit `BEGIN/COMMIT` commands included in the query string to divide it into multiple transactions. (See [Section 56.2.2.1](#) for more details about how the server handles multi-query strings.) Note however that the returned `PGresult` structure describes only the result of the last command executed from the string. Should one of the commands fail, processing of the string stops with it and the returned `PGresult` describes the error condition.

[PQexecParams](#)

Submits a command to the server and waits for the result, with the ability to pass parameters separately from the SQL command text.

```
PGresult *PQexecParams(PGconn *conn,
                        const char *command,
                        int nParams,
                        const Oid *paramTypes,
                        const char * const *paramValues,
                        const int *paramLengths,
                        const int *paramFormats,
                        int resultFormat);
```

[PQexecParams](#) is like [PQexec](#), but offers additional functionality: parameter values can be specified separately from the command string proper, and query results can be requested in either text or binary format.

The function arguments are:

conn

The connection object to send the command through.

command

The SQL command string to be executed. If parameters are used, they are referred to in the command string as `$1`, `$2`, etc.

nParams

The number of parameters supplied; it is the length of the arrays `paramTypes[]`, `paramValues[]`, `paramLengths[]`, and `paramFormats[]`. (The array pointers can be `NULL` when `nParams` is zero.)

paramTypes[]

Specifies, by OID, the data types to be assigned to the parameter symbols. If *paramTypes* is NULL, or any particular element in the array is zero, the server infers a data type for the parameter symbol in the same way it would do for an untyped literal string.

paramValues[]

Specifies the actual values of the parameters. A null pointer in this array means the corresponding parameter is null; otherwise the pointer points to a zero-terminated text string (for text format) or binary data in the format expected by the server (for binary format).

paramLengths[]

Specifies the actual data lengths of binary-format parameters. It is ignored for null parameters and text-format parameters. The array pointer can be null when there are no binary parameters.

paramFormats[]

Specifies whether parameters are text (put a zero in the array entry for the corresponding parameter) or binary (put a one in the array entry for the corresponding parameter). If the array pointer is null then all parameters are presumed to be text strings.

Values passed in binary format require knowledge of the internal representation expected by the backend. For example, integers must be passed in network byte order. Passing `numeric` values requires knowledge of the server storage format.

resultFormat

Specify zero to obtain results in text format, or one to obtain results in binary format. (There is not currently a provision to obtain different result columns in different formats, although that is possible in the underlying protocol.)

The primary advantage of `PQexecParams` over `PQexec` is that parameter values can be separated from the command string, thus avoiding the need for tedious and error-prone quoting and escaping.

Unlike `PQexec`, `PQexecParams` allows at most one SQL command in the given string. (There can be semicolons in it, but not more than one nonempty command.) This is a limitation of the underlying protocol, but has some usefulness as an extra defense against SQL-injection attacks.

Tip

Specifying parameter types via OIDs is tedious, particularly if you prefer not to hard-wire particular OID values into your program. However, you can avoid doing so even in cases where the server by itself cannot determine the type of the parameter, or chooses a different type than you want. In the SQL command text, attach an explicit cast to the parameter symbol to show what data type you will send. For example:

```
SELECT * FROM mytable WHERE x = $1::bigint;
```

This forces parameter `$1` to be treated as `bigint`, whereas by default it would be assigned the same type as `x`. Forcing the parameter type decision, either this way or by specifying a numeric type OID, is strongly recommended when sending parameter values in binary format, because binary format has less redundancy than text format and so there is less chance that the server will detect a type mismatch mistake for you.

`PQprepare`

Submits a request to create a prepared statement with the given parameters, and waits for completion.

```
PGresult *PQprepare(PGconn *conn,
                    const char *stmtName,
```

```
const char *query,  
int nParams,  
const Oid *paramTypes);
```

[PQprepare](#) creates a prepared statement for later execution with [PQexecPrepared](#). This feature allows commands to be executed repeatedly without being parsed and planned each time; see [PREPARE](#) for details.

The function creates a prepared statement named *stmtName* from the *query* string, which must contain a single SQL command. *stmtName* can be "" to create an unnamed statement, in which case any pre-existing unnamed statement is automatically replaced; otherwise it is an error if the statement name is already defined in the current session. If any parameters are used, they are referred to in the query as \$1, \$2, etc. *nParams* is the number of parameters for which types are pre-specified in the array *paramTypes[]*. (The array pointer can be NULL when *nParams* is zero.) *paramTypes[]* specifies, by OID, the data types to be assigned to the parameter symbols. If *paramTypes* is NULL, or any particular element in the array is zero, the server assigns a data type to the parameter symbol in the same way it would do for an untyped literal string. Also, the query can use parameter symbols with numbers higher than *nParams*; data types will be inferred for these symbols as well. (See [PQdescribePrepared](#) for a means to find out what data types were inferred.)

As with [PQexec](#), the result is normally a `PGresult` object whose contents indicate server-side success or failure. A null result indicates out-of-memory or inability to send the command at all. Use [PQerrorMessage](#) to get more information about such errors.

Prepared statements for use with [PQexecPrepared](#) can also be created by executing SQL [PREPARE](#) statements.

`PQexecPrepared`

Sends a request to execute a prepared statement with given parameters, and waits for the result.

```
PGresult *PQexecPrepared(PGconn *conn,  
                          const char *stmtName,  
                          int nParams,  
                          const char * const *paramValues,  
                          const int *paramLengths,  
                          const int *paramFormats,  
                          int resultFormat);
```

[PQexecPrepared](#) is like [PQexecParams](#), but the command to be executed is specified by naming a previously-prepared statement, instead of giving a query string. This feature allows commands that will be used repeatedly to be parsed and planned just once, rather than each time they are executed. The statement must have been prepared previously in the current session.

The parameters are identical to [PQexecParams](#), except that the name of a prepared statement is given instead of a query string, and the *paramTypes[]* parameter is not present (it is not needed since the prepared statement's parameter types were determined when it was created).

`PQdescribePrepared`

Submits a request to obtain information about the specified prepared statement, and waits for completion.

```
PGresult *PQdescribePrepared(PGconn *conn, const char *stmtName);
```

[PQdescribePrepared](#) allows an application to obtain information about a previously prepared statement.

stmtName can be "" or NULL to reference the unnamed statement, otherwise it must be the name of an existing prepared statement. On success, a `PGresult` with status `PGRES_COMMAND_OK` is returned. The functions [PQnparams](#) and [PQparamtype](#) can be applied to this `PGresult` to obtain information

about the parameters of the prepared statement, and the functions `PQnfields`, `PQfname`, `PQftype`, etc. provide information about the result columns (if any) of the statement.

`PQdescribePortal`

Submits a request to obtain information about the specified portal, and waits for completion.

```
PGresult *PQdescribePortal(PGconn *conn, const char *portalName);
```

`PQdescribePortal` allows an application to obtain information about a previously created portal. (libpq does not provide any direct access to portals, but you can use this function to inspect the properties of a cursor created with a `DECLARE CURSOR SQL` command.)

`portalName` can be "" or `NULL` to reference the unnamed portal, otherwise it must be the name of an existing portal. On success, a `PGresult` with status `PGRES_COMMAND_OK` is returned. The functions `PQnfields`, `PQfname`, `PQftype`, etc. can be applied to the `PGresult` to obtain information about the result columns (if any) of the portal.

`PQclosePrepared`

Submits a request to close the specified prepared statement, and waits for completion.

```
PGresult *PQclosePrepared(PGconn *conn, const char *stmtName);
```

`PQclosePrepared` allows an application to close a previously prepared statement. Closing a statement releases all of its associated resources on the server and allows its name to be reused.

`stmtName` can be "" or `NULL` to reference the unnamed statement. It is fine if no statement exists with this name, in that case the operation is a no-op. On success, a `PGresult` with status `PGRES_COMMAND_OK` is returned.

`PQclosePortal`

Submits a request to close the specified portal, and waits for completion.

```
PGresult *PQclosePortal(PGconn *conn, const char *portalName);
```

`PQclosePortal` allows an application to trigger a close of a previously created portal. Closing a portal releases all of its associated resources on the server and allows its name to be reused. (libpq does not provide any direct access to portals, but you can use this function to close a cursor created with a `DECLARE CURSOR SQL` command.)

`portalName` can be "" or `NULL` to reference the unnamed portal. It is fine if no portal exists with this name, in that case the operation is a no-op. On success, a `PGresult` with status `PGRES_COMMAND_OK` is returned.

The `PGresult` structure encapsulates the result returned by the server. libpq application programmers should be careful to maintain the `PGresult` abstraction. Use the accessor functions below to get at the contents of `PGresult`. Avoid directly referencing the fields of the `PGresult` structure because they are subject to change in the future.

`PQresultStatus`

Returns the result status of the command.

```
ExecStatusType PQresultStatus(const PGresult *res);
```

`PQresultStatus` can return one of the following values:

`PGRES_EMPTY_QUERY`

The string sent to the server was empty.

`PGRES_COMMAND_OK`

Successful completion of a command returning no data.

`PGRES_TUPLES_OK`

Successful completion of a command returning data (such as a `SELECT` or `SHOW`).

`PGRES_COPY_OUT`

Copy Out (from server) data transfer started.

`PGRES_COPY_IN`

Copy In (to server) data transfer started.

`PGRES_BAD_RESPONSE`

The server's response was not understood.

`PGRES_NONFATAL_ERROR`

A nonfatal error (a notice or warning) occurred.

`PGRES_FATAL_ERROR`

A fatal error occurred.

`PGRES_COPY_BOTH`

Copy In/Out (to and from server) data transfer started. This feature is currently used only for streaming replication, so this status should not occur in ordinary applications.

`PGRES_SINGLE_TUPLE`

The `PGresult` contains a single result tuple from the current command. This status occurs only when single-row mode has been selected for the query (see [Section 35.6](#)).

`PGRES_TUPLES_CHUNK`

The `PGresult` contains several result tuples from the current command. This status occurs only when chunked mode has been selected for the query (see [Section 35.6](#)). The number of tuples will not exceed the limit passed to `PQsetChunkedRowsMode`.

`PGRES_PIPELINE_SYNC`

The `PGresult` represents a synchronization point in pipeline mode, requested by either `PQpipelineSync` or `PQsendPipelineSync`. This status occurs only when pipeline mode has been selected.

`PGRES_PIPELINE_ABORTED`

The `PGresult` represents a pipeline that has received an error from the server. `PQgetResult` must be called repeatedly, and each time it will return this status code until the end of the current pipeline, at which point it will return `PGRES_PIPELINE_SYNC` and normal processing can resume.

If the result status is `PGRES_TUPLES_OK`, `PGRES_SINGLE_TUPLE`, or `PGRES_TUPLES_CHUNK`, then the functions described below can be used to retrieve the rows returned by the query. Note that a `SELECT` command that happens to retrieve zero rows still shows `PGRES_TUPLES_OK`. `PGRES_COMMAND_OK` is for commands that can never return rows (`INSERT` or `UPDATE` without a `RETURNING` clause, etc.). A response of `PGRES_EMPTY_QUERY` might indicate a bug in the client software.

A result of status `PGRES_NONFATAL_ERROR` will never be returned directly by `PQexec` or other query execution functions; results of this kind are instead passed to the notice processor (see [Section 35.13](#)).

`PQresStatus`

Converts the enumerated type returned by `PQresultStatus` into a string constant describing the status code. The caller should not free the result.

```
char *PQresStatus(ExecStatusType status);
```

PQresultErrorMessage

Returns the error message associated with the command, or an empty string if there was no error.

```
char *PQresultErrorMessage(const PGresult *res);
```

If there was an error, the returned string will include a trailing newline. The caller should not free the result directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

Immediately following a `PQexec` or `PQgetResult` call, `PQerrorMessage` (on the connection) will return the same string as `PQresultErrorMessage` (on the result). However, a `PGresult` will retain its error message until destroyed, whereas the connection's error message will change when subsequent operations are done. Use `PQresultErrorMessage` when you want to know the status associated with a particular `PGresult`; use `PQerrorMessage` when you want to know the status from the latest operation on the connection.

PQresultVerboseErrorMessage

Returns a reformatted version of the error message associated with a `PGresult` object.

```
char *PQresultVerboseErrorMessage(const PGresult *res,
                                  PGVerbosity verbosity,
                                  PGContextVisibility show_context);
```

In some situations a client might wish to obtain a more detailed version of a previously-reported error. `PQresultVerboseErrorMessage` addresses this need by computing the message that would have been produced by `PQresultErrorMessage` if the specified verbosity settings had been in effect for the connection when the given `PGresult` was generated. If the `PGresult` is not an error result, "PGresult is not an error result" is reported instead. The returned string includes a trailing newline.

Unlike most other functions for extracting data from a `PGresult`, the result of this function is a freshly allocated string. The caller must free it using `PQfreemem()` when the string is no longer needed.

A NULL return is possible if there is insufficient memory.

PQresultErrorField

Returns an individual field of an error report.

```
char *PQresultErrorField(const PGresult *res, int fieldcode);
```

fieldcode is an error field identifier; see the symbols listed below. NULL is returned if the `PGresult` is not an error or warning result, or does not include the specified field. Field values will normally not include a trailing newline. The caller should not free the result directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

The following field codes are available:

PG_DIAG_SEVERITY

The severity; the field contents are `ERROR`, `FATAL`, or `PANIC` (in an error message), or `WARNING`, `NOTICE`, `DEBUG`, `INFO`, or `LOG` (in a notice message), or a localized translation of one of these. Always present.

PG_DIAG_SEVERITY_NONLOCALIZED

The severity; the field contents are `ERROR`, `FATAL`, or `PANIC` (in an error message), or `WARNING`, `NOTICE`, `DEBUG`, `INFO`, or `LOG` (in a notice message). This is identical to the `PG_DIAG_SEVERITY` field except that the contents are never localized. This is present only in reports generated by Postgres Pro versions 9.6 and later.

PG_DIAG_SQLSTATE

The `SQLSTATE` code for the error. The `SQLSTATE` code identifies the type of error that has occurred; it can be used by front-end applications to perform specific operations (such as error

handling) in response to a particular database error. For a list of the possible SQLSTATE codes, see [Appendix A](#). This field is not localizable, and is always present.

PG_DIAG_MESSAGE_PRIMARY

The primary human-readable error message (typically one line). Always present.

PG_DIAG_MESSAGE_DETAIL

Detail: an optional secondary error message carrying more detail about the problem. Might run to multiple lines.

PG_DIAG_MESSAGE_HINT

Hint: an optional suggestion what to do about the problem. This is intended to differ from detail in that it offers advice (potentially inappropriate) rather than hard facts. Might run to multiple lines.

PG_DIAG_STATEMENT_POSITION

A string containing a decimal integer indicating an error cursor position as an index into the original statement string. The first character has index 1, and positions are measured in characters not bytes.

PG_DIAG_INTERNAL_POSITION

This is defined the same as the `PG_DIAG_STATEMENT_POSITION` field, but it is used when the cursor position refers to an internally generated command rather than the one submitted by the client. The `PG_DIAG_INTERNAL_QUERY` field will always appear when this field appears.

PG_DIAG_INTERNAL_QUERY

The text of a failed internally-generated command. This could be, for example, an SQL query issued by a PL/pgSQL function.

PG_DIAG_CONTEXT

An indication of the context in which the error occurred. Presently this includes a call stack traceback of active procedural language functions and internally-generated queries. The trace is one entry per line, most recent first.

PG_DIAG_SCHEMA_NAME

If the error was associated with a specific database object, the name of the schema containing that object, if any.

PG_DIAG_TABLE_NAME

If the error was associated with a specific table, the name of the table. (Refer to the schema name field for the name of the table's schema.)

PG_DIAG_COLUMN_NAME

If the error was associated with a specific table column, the name of the column. (Refer to the schema and table name fields to identify the table.)

PG_DIAG_DATATYPE_NAME

If the error was associated with a specific data type, the name of the data type. (Refer to the schema name field for the name of the data type's schema.)

PG_DIAG_CONSTRAINT_NAME

If the error was associated with a specific constraint, the name of the constraint. Refer to fields listed above for the associated table or domain. (For this purpose, indexes are treated as constraints, even if they weren't created with constraint syntax.)

`PG_DIAG_SOURCE_FILE`

The file name of the source-code location where the error was reported.

`PG_DIAG_SOURCE_LINE`

The line number of the source-code location where the error was reported.

`PG_DIAG_SOURCE_FUNCTION`

The name of the source-code function reporting the error.

Note

The fields for schema name, table name, column name, data type name, and constraint name are supplied only for a limited number of error types; see [Appendix A](#). Do not assume that the presence of any of these fields guarantees the presence of another field. Core error sources observe the interrelationships noted above, but user-defined functions may use these fields in other ways. In the same vein, do not assume that these fields denote contemporary objects in the current database.

The client is responsible for formatting displayed information to meet its needs; in particular it should break long lines as needed. Newline characters appearing in the error message fields should be treated as paragraph breaks, not line breaks.

Errors generated internally by libpq will have severity and primary message, but typically no other fields.

Note that error fields are only available from `PGresult` objects, not `PGconn` objects; there is no `PQerrorField` function.

`PQclear`

Frees the storage associated with a `PGresult`. Every command result should be freed via `PQclear` when it is no longer needed.

```
void PQclear(PGresult *res);
```

If the argument is a `NULL` pointer, no operation is performed.

You can keep a `PGresult` object around for as long as you need it; it does not go away when you issue a new command, nor even if you close the connection. To get rid of it, you must call `PQclear`. Failure to do this will result in memory leaks in your application.

35.3.2. Retrieving Query Result Information

These functions are used to extract information from a `PGresult` object that represents a successful query result (that is, one that has status `PGRES_TUPLES_OK`, `PGRES_SINGLE_TUPLE`, or `PGRES_TUPLES_CHUNK`). They can also be used to extract information from a successful Describe operation: a Describe's result has all the same column information that actual execution of the query would provide, but it has zero rows. For objects with other status values, these functions will act as though the result has zero rows and zero columns.

`PQntuples`

Returns the number of rows (tuples) in the query result. (Note that `PGresult` objects are limited to no more than `INT_MAX` rows, so an `int` result is sufficient.)

```
int PQntuples(const PGresult *res);
```

`PQnfields`

Returns the number of columns (fields) in each row of the query result.

```
int PQnfields(const PGresult *res);
```

PQfname

Returns the column name associated with the given column number. Column numbers start at 0. The caller should not free the result directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

```
char *PQfname(const PGresult *res,
              int column_number);
```

NULL is returned if the column number is out of range.

PQfnumber

Returns the column number associated with the given column name.

```
int PQfnumber(const PGresult *res,
              const char *column_name);
```

-1 is returned if the given name does not match any column.

The given name is treated like an identifier in an SQL command, that is, it is downcased unless double-quoted. For example, given a query result generated from the SQL command:

```
SELECT 1 AS FOO, 2 AS "BAR";
```

we would have the results:

<code>PQfname(res, 0)</code>	<code>foo</code>
<code>PQfname(res, 1)</code>	<code>BAR</code>
<code>PQfnumber(res, "FOO")</code>	<code>0</code>
<code>PQfnumber(res, "foo")</code>	<code>0</code>
<code>PQfnumber(res, "BAR")</code>	<code>-1</code>
<code>PQfnumber(res, "\"BAR\"")</code>	<code>1</code>

PQftable

Returns the OID of the table from which the given column was fetched. Column numbers start at 0.

```
Oid PQftable(const PGresult *res,
             int column_number);
```

`InvalidOid` is returned if the column number is out of range, or if the specified column is not a simple reference to a table column. You can query the system table `pg_class` to determine exactly which table is referenced.

The type `Oid` and the constant `InvalidOid` will be defined when you include the libpq header file. They will both be some integer type.

PQftablecol

Returns the column number (within its table) of the column making up the specified query result column. Query-result column numbers start at 0, but table columns have nonzero numbers.

```
int PQftablecol(const PGresult *res,
                int column_number);
```

Zero is returned if the column number is out of range, or if the specified column is not a simple reference to a table column.

PQfformat

Returns the format code indicating the format of the given column. Column numbers start at 0.

```
int PQfformat(const PGresult *res,
```

```
int column_number);
```

Format code zero indicates textual data representation, while format code one indicates binary representation. (Other codes are reserved for future definition.)

PQftype

Returns the data type associated with the given column number. The integer returned is the internal OID number of the type. Column numbers start at 0.

```
Oid PQftype(const PGresult *res,
            int column_number);
```

You can query the system table `pg_type` to obtain the names and properties of the various data types.

PQfmod

Returns the type modifier of the column associated with the given column number. Column numbers start at 0.

```
int PQfmod(const PGresult *res,
           int column_number);
```

The interpretation of modifier values is type-specific; they typically indicate precision or size limits. The value -1 is used to indicate “no information available”. Most data types do not use modifiers, in which case the value is always -1.

PQfsize

Returns the size in bytes of the column associated with the given column number. Column numbers start at 0.

```
int PQfsize(const PGresult *res,
            int column_number);
```

`PQfsize` returns the space allocated for this column in a database row, in other words the size of the server's internal representation of the data type. (Accordingly, it is not really very useful to clients.) A negative value indicates the data type is variable-length.

PQbinaryTuples

Returns 1 if the `PGresult` contains binary data and 0 if it contains text data.

```
int PQbinaryTuples(const PGresult *res);
```

This function is deprecated (except for its use in connection with `COPY`), because it is possible for a single `PGresult` to contain text data in some columns and binary data in others. `PQfformat` is preferred. `PQbinaryTuples` returns 1 only if all columns of the result are binary (format 1).

PQgetvalue

Returns a single field value of one row of a `PGresult`. Row and column numbers start at 0. The caller should not free the result directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

```
char *PQgetvalue(const PGresult *res,
                 int row_number,
                 int column_number);
```

For data in text format, the value returned by `PQgetvalue` is a null-terminated character string representation of the field value. For data in binary format, the value is in the binary representation determined by the data type's `typsend` and `typreceive` functions. (The value is actually followed by a zero byte in this case too, but that is not ordinarily useful, since the value is likely to contain embedded nulls.)

An empty string is returned if the field value is null. See [PQgetisnull](#) to distinguish null values from empty-string values.

The pointer returned by [PQgetvalue](#) points to storage that is part of the `PGresult` structure. One should not modify the data it points to, and one must explicitly copy the data into other storage if it is to be used past the lifetime of the `PGresult` structure itself.

`PQgetisnull`

Tests a field for a null value. Row and column numbers start at 0.

```
int PQgetisnull(const PGresult *res,
                int row_number,
                int column_number);
```

This function returns 1 if the field is null and 0 if it contains a non-null value. (Note that [PQgetvalue](#) will return an empty string, not a null pointer, for a null field.)

`PQgetlength`

Returns the actual length of a field value in bytes. Row and column numbers start at 0.

```
int PQgetlength(const PGresult *res,
                int row_number,
                int column_number);
```

This is the actual data length for the particular data value, that is, the size of the object pointed to by [PQgetvalue](#). For text data format this is the same as `strlen()`. For binary format this is essential information. Note that one should *not* rely on [PQfsize](#) to obtain the actual data length.

`PQnparams`

Returns the number of parameters of a prepared statement.

```
int PQnparams(const PGresult *res);
```

This function is only useful when inspecting the result of [PQdescribePrepared](#). For other types of results it will return zero.

`PQparamtype`

Returns the data type of the indicated statement parameter. Parameter numbers start at 0.

```
Oid PQparamtype(const PGresult *res, int param_number);
```

This function is only useful when inspecting the result of [PQdescribePrepared](#). For other types of results it will return zero.

`PQprint`

Prints out all the rows and, optionally, the column names to the specified output stream.

```
void PQprint(FILE *fout,          /* output stream */
             const PGresult *res,
             const PQprintOpt *po);

typedef struct
{
    pqbool header;          /* print output field headings and row count */
    pqbool align;           /* fill align the fields */
    pqbool standard;        /* old brain dead format */
    pqbool html3;           /* output HTML tables */
    pqbool expanded;        /* expand tables */
    pqbool pager;           /* use pager for output if needed */
    char *fieldSep;         /* field separator */
}
```



```
char    *tableOpt;    /* attributes for HTML table element */
char    *caption;     /* HTML table caption */
char    **fieldName; /* null-terminated array of replacement field names */
} PQprintOpt;
```

This function was formerly used by `psql` to print query results, but this is no longer the case. Note that it assumes all the data is in text format.

35.3.3. Retrieving Other Result Information

These functions are used to extract other information from `PGresult` objects.

`PQcmdStatus`

Returns the command status tag from the SQL command that generated the `PGresult`.

```
char *PQcmdStatus(PGresult *res);
```

Commonly this is just the name of the command, but it might include additional data such as the number of rows processed. The caller should not free the result directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

`PQcmdTuples`

Returns the number of rows affected by the SQL command.

```
char *PQcmdTuples(PGresult *res);
```

This function returns a string containing the number of rows affected by the SQL statement that generated the `PGresult`. This function can only be used following the execution of a `SELECT`, `CREATE TABLE AS`, `INSERT`, `UPDATE`, `DELETE`, `MERGE`, `MOVE`, `FETCH`, or `COPY` statement, or an `EXECUTE` of a prepared query that contains an `INSERT`, `UPDATE`, `DELETE`, or `MERGE` statement. If the command that generated the `PGresult` was anything else, `PQcmdTuples` returns an empty string. The caller should not free the return value directly. It will be freed when the associated `PGresult` handle is passed to `PQclear`.

`PQoidValue`

Returns the OID of the inserted row, if the SQL command was an `INSERT` that inserted exactly one row into a table that has OIDs, or a `EXECUTE` of a prepared query containing a suitable `INSERT` statement. Otherwise, this function returns `InvalidOid`. This function will also return `InvalidOid` if the table affected by the `INSERT` statement does not contain OIDs.

```
Oid PQoidValue(const PGresult *res);
```

`PQoidStatus`

This function is deprecated in favor of `PQoidValue` and is not thread-safe. It returns a string with the OID of the inserted row, while `PQoidValue` returns the OID value.

```
char *PQoidStatus(const PGresult *res);
```

35.3.4. Escaping Strings for Inclusion in SQL Commands

`PQescapeLiteral`

```
char *PQescapeLiteral(PGconn *conn, const char *str, size_t length);
```

`PQescapeLiteral` escapes a string for use within an SQL command. This is useful when inserting data values as literal constants in SQL commands. Certain characters (such as quotes and backslashes) must be escaped to prevent them from being interpreted specially by the SQL parser. `PQescapeLiteral` performs this operation.

`PQescapeLiteral` returns an escaped version of the `str` parameter in memory allocated with `malloc()`. This memory should be freed using `PQfreemem()` when the result is no longer needed. A ter-

minating zero byte is not required, and should not be counted in *length*. (If a terminating zero byte is found before *length* bytes are processed, `PQescapeLiteral` stops at the zero; the behavior is thus rather like `strncpy`.) The return string has all special characters replaced so that they can be properly processed by the Postgres Pro string literal parser. A terminating zero byte is also added. The single quotes that must surround Postgres Pro string literals are included in the result string.

On error, `PQescapeLiteral` returns `NULL` and a suitable message is stored in the *conn* object.

Tip

It is especially important to do proper escaping when handling strings that were received from an untrustworthy source. Otherwise there is a security risk: you are vulnerable to “SQL injection” attacks wherein unwanted SQL commands are fed to your database.

Note that it is neither necessary nor correct to do escaping when a data value is passed as a separate parameter in `PQexecParams` or its sibling routines.

`PQescapeIdentifier`

```
char *PQescapeIdentifier(PGconn *conn, const char *str, size_t length);
```

`PQescapeIdentifier` escapes a string for use as an SQL identifier, such as a table, column, or function name. This is useful when a user-supplied identifier might contain special characters that would otherwise not be interpreted as part of the identifier by the SQL parser, or when the identifier might contain upper case characters whose case should be preserved.

`PQescapeIdentifier` returns a version of the *str* parameter escaped as an SQL identifier in memory allocated with `malloc()`. This memory must be freed using `PQfreemem()` when the result is no longer needed. A terminating zero byte is not required, and should not be counted in *length*. (If a terminating zero byte is found before *length* bytes are processed, `PQescapeIdentifier` stops at the zero; the behavior is thus rather like `strncpy`.) The return string has all special characters replaced so that it will be properly processed as an SQL identifier. A terminating zero byte is also added. The return string will also be surrounded by double quotes.

On error, `PQescapeIdentifier` returns `NULL` and a suitable message is stored in the *conn* object.

Tip

As with string literals, to prevent SQL injection attacks, SQL identifiers must be escaped when they are received from an untrustworthy source.

`PQescapeStringConn`

```
size_t PQescapeStringConn(PGconn *conn,
                           char *to, const char *from, size_t length,
                           int *error);
```

`PQescapeStringConn` escapes string literals, much like `PQescapeLiteral`. Unlike `PQescapeLiteral`, the caller is responsible for providing an appropriately sized buffer. Furthermore, `PQescapeStringConn` does not generate the single quotes that must surround Postgres Pro string literals; they should be provided in the SQL command that the result is inserted into. The parameter *from* points to the first character of the string that is to be escaped, and the *length* parameter gives the number of bytes in this string. A terminating zero byte is not required, and should not be counted in *length*. (If a terminating zero byte is found before *length* bytes are processed, `PQescapeStringConn` stops at the zero; the behavior is thus rather like `strncpy`.) *to* shall point to a buffer that is able to hold at least one more byte than twice the value of *length*, otherwise the behavior is undefined. Behavior is likewise undefined if the *to* and *from* strings overlap.

If the `error` parameter is not `NULL`, then `*error` is set to zero on success, nonzero on error. Presently the only possible error conditions involve invalid multibyte encoding in the source string. The output string is still generated on error, but it can be expected that the server will reject it as malformed. On error, a suitable message is stored in the `conn` object, whether or not `error` is `NULL`.

`PQescapeStringConn` returns the number of bytes written to `to`, not including the terminating zero byte.

`PQescapeString`

`PQescapeString` is an older, deprecated version of `PQescapeStringConn`.

```
size_t PQescapeString (char *to, const char *from, size_t length);
```

The only difference from `PQescapeStringConn` is that `PQescapeString` does not take `PGconn` or `error` parameters. Because of this, it cannot adjust its behavior depending on the connection properties (such as character encoding) and therefore *it might give the wrong results*. Also, it has no way to report error conditions.

`PQescapeString` can be used safely in client programs that work with only one Postgres Pro connection at a time (in this case it can find out what it needs to know “behind the scenes”). In other contexts it is a security hazard and should be avoided in favor of `PQescapeStringConn`.

`PQescapeByteaConn`

Escapes binary data for use within an SQL command with the type `bytea`. As with `PQescapeStringConn`, this is only used when inserting data directly into an SQL command string.

```
unsigned char *PQescapeByteaConn (PGconn *conn,
                                   const unsigned char *from,
                                   size_t from_length,
                                   size_t *to_length);
```

Certain byte values must be escaped when used as part of a `bytea` literal in an SQL statement. `PQescapeByteaConn` escapes bytes using either hex encoding or backslash escaping. See [Section 8.4](#) for more information.

The `from` parameter points to the first byte of the string that is to be escaped, and the `from_length` parameter gives the number of bytes in this binary string. (A terminating zero byte is neither necessary nor counted.) The `to_length` parameter points to a variable that will hold the resultant escaped string length. This result string length includes the terminating zero byte of the result.

`PQescapeByteaConn` returns an escaped version of the `from` parameter binary string in memory allocated with `malloc()`. This memory should be freed using `PQfreemem()` when the result is no longer needed. The return string has all special characters replaced so that they can be properly processed by the Postgres Pro string literal parser, and the `bytea` input function. A terminating zero byte is also added. The single quotes that must surround Postgres Pro string literals are not part of the result string.

On error, a null pointer is returned, and a suitable error message is stored in the `conn` object. Currently, the only possible error is insufficient memory for the result string.

`PQescapeBytea`

`PQescapeBytea` is an older, deprecated version of `PQescapeByteaConn`.

```
unsigned char *PQescapeBytea (const unsigned char *from,
                              size_t from_length,
                              size_t *to_length);
```

The only difference from `PQescapeByteaConn` is that `PQescapeBytea` does not take a `PGconn` parameter. Because of this, `PQescapeBytea` can only be used safely in client programs that use a single

Postgres Pro connection at a time (in this case it can find out what it needs to know “behind the scenes”). It *might give the wrong results* if used in programs that use multiple database connections (use `PQescapeByteaConn` in such cases).

`PQunescapeBytea`

Converts a string representation of binary data into binary data — the reverse of `PQescapeBytea`. This is needed when retrieving `bytea` data in text format, but not when retrieving it in binary format.

```
unsigned char *PQunescapeBytea(const unsigned char *from, size_t *to_length);
```

The `from` parameter points to a string such as might be returned by `PQgetvalue` when applied to a `bytea` column. `PQunescapeBytea` converts this string representation into its binary representation. It returns a pointer to a buffer allocated with `malloc()`, or `NULL` on error, and puts the size of the buffer in `to_length`. The result must be freed using `PQfreemem` when it is no longer needed.

This conversion is not exactly the inverse of `PQescapeBytea`, because the string is not expected to be “escaped” when received from `PQgetvalue`. In particular this means there is no need for string quoting considerations, and so no need for a `PGconn` parameter.

35.4. Asynchronous Command Processing

The `PQexec` function is adequate for submitting commands in normal, synchronous applications. It has a few deficiencies, however, that can be of importance to some users:

- `PQexec` waits for the command to be completed. The application might have other work to do (such as maintaining a user interface), in which case it won't want to block waiting for the response.
- Since the execution of the client application is suspended while it waits for the result, it is hard for the application to decide that it would like to try to cancel the ongoing command. (It can be done from a signal handler, but not otherwise.)
- `PQexec` can return only one `PGresult` structure. If the submitted command string contains multiple SQL commands, all but the last `PGresult` are discarded by `PQexec`.
- `PQexec` always collects the command's entire result, buffering it in a single `PGresult`. While this simplifies error-handling logic for the application, it can be impractical for results containing many rows.

Applications that do not like these limitations can instead use the underlying functions that `PQexec` is built from: `PQsendQuery` and `PQgetResult`. There are also `PQsendQueryParams`, `PQsendPrepare`, `PQsendQueryPrepared`, `PQsendDescribePrepared`, `PQsendDescribePortal`, `PQsendClosePrepared`, and `PQsendClosePortal`, which can be used with `PQgetResult` to duplicate the functionality of `PQexecParams`, `PQprepare`, `PQexecPrepared`, `PQdescribePrepared`, `PQdescribePortal`, `PQclosePrepared`, and `PQclosePortal` respectively.

`PQsendQuery`

Submits a command to the server without waiting for the result(s). 1 is returned if the command was successfully dispatched and 0 if not (in which case, use `PQerrorMessage` to get more information about the failure).

```
int PQsendQuery(PGconn *conn, const char *command);
```

After successfully calling `PQsendQuery`, call `PQgetResult` one or more times to obtain the results. `PQsendQuery` cannot be called again (on the same connection) until `PQgetResult` has returned a null pointer, indicating that the command is done.

In pipeline mode, this function is disallowed.

`PQsendQueryParams`

Submits a command and separate parameters to the server without waiting for the result(s).

```
int PQsendQueryParams(PGconn *conn,
                     const char *command,
                     int nParams,
                     const Oid *paramTypes,
                     const char * const *paramValues,
                     const int *paramLengths,
                     const int *paramFormats,
                     int resultFormat);
```

This is equivalent to [PQsendQuery](#) except that query parameters can be specified separately from the query string. The function's parameters are handled identically to [PQexecParams](#). Like [PQexecParams](#), it allows only one command in the query string.

[PQsendPrepare](#)

Sends a request to create a prepared statement with the given parameters, without waiting for completion.

```
int PQsendPrepare(PGconn *conn,
                 const char *stmtName,
                 const char *query,
                 int nParams,
                 const Oid *paramTypes);
```

This is an asynchronous version of [PQprepare](#): it returns 1 if it was able to dispatch the request, and 0 if not. After a successful call, call [PQgetResult](#) to determine whether the server successfully created the prepared statement. The function's parameters are handled identically to [PQprepare](#).

[PQsendQueryPrepared](#)

Sends a request to execute a prepared statement with given parameters, without waiting for the result(s).

```
int PQsendQueryPrepared(PGconn *conn,
                       const char *stmtName,
                       int nParams,
                       const char * const *paramValues,
                       const int *paramLengths,
                       const int *paramFormats,
                       int resultFormat);
```

This is similar to [PQsendQueryParams](#), but the command to be executed is specified by naming a previously-prepared statement, instead of giving a query string. The function's parameters are handled identically to [PQexecPrepared](#).

[PQsendDescribePrepared](#)

Submits a request to obtain information about the specified prepared statement, without waiting for completion.

```
int PQsendDescribePrepared(PGconn *conn, const char *stmtName);
```

This is an asynchronous version of [PQdescribePrepared](#): it returns 1 if it was able to dispatch the request, and 0 if not. After a successful call, call [PQgetResult](#) to obtain the results. The function's parameters are handled identically to [PQdescribePrepared](#).

[PQsendDescribePortal](#)

Submits a request to obtain information about the specified portal, without waiting for completion.

```
int PQsendDescribePortal(PGconn *conn, const char *portalName);
```

This is an asynchronous version of [PQdescribePortal](#): it returns 1 if it was able to dispatch the request, and 0 if not. After a successful call, call [PQgetResult](#) to obtain the results. The function's parameters are handled identically to [PQdescribePortal](#).

PQsendClosePrepared

Submits a request to close the specified prepared statement, without waiting for completion.

```
int PQsendClosePrepared(PGconn *conn, const char *stmtName);
```

This is an asynchronous version of [PQclosePrepared](#): it returns 1 if it was able to dispatch the request, and 0 if not. After a successful call, call [PQgetResult](#) to obtain the results. The function's parameters are handled identically to [PQclosePrepared](#).

PQsendClosePortal

Submits a request to close specified portal, without waiting for completion.

```
int PQsendClosePortal(PGconn *conn, const char *portalName);
```

This is an asynchronous version of [PQclosePortal](#): it returns 1 if it was able to dispatch the request, and 0 if not. After a successful call, call [PQgetResult](#) to obtain the results. The function's parameters are handled identically to [PQclosePortal](#).

PQgetResult

Waits for the next result from a prior [PQsendQuery](#), [PQsendQueryParams](#), [PQsendPrepare](#), [PQsendQueryPrepared](#), [PQsendDescribePrepared](#), [PQsendDescribePortal](#), [PQsendClosePrepared](#), [PQsendClosePortal](#), [PQsendPipelineSync](#), or [PQpipelineSync](#) call, and returns it. A null pointer is returned when the command is complete and there will be no more results.

```
PGresult *PQgetResult(PGconn *conn);
```

[PQgetResult](#) must be called repeatedly until it returns a null pointer, indicating that the command is done. (If called when no command is active, [PQgetResult](#) will just return a null pointer at once.) Each non-null result from [PQgetResult](#) should be processed using the same `PGresult` accessor functions previously described. Don't forget to free each result object with [PQclear](#) when done with it. Note that [PQgetResult](#) will block only if a command is active and the necessary response data has not yet been read by [PQconsumeInput](#).

In pipeline mode, [PQgetResult](#) will return normally unless an error occurs; for any subsequent query sent after the one that caused the error until (and excluding) the next synchronization point, a special result of type `PGRES_PIPELINE_ABORTED` will be returned, and a null pointer will be returned after it. When the pipeline synchronization point is reached, a result of type `PGRES_PIPELINE_SYNC` will be returned. The result of the next query after the synchronization point follows immediately (that is, no null pointer is returned after the synchronization point).

Note

Even when [PQresultStatus](#) indicates a fatal error, [PQgetResult](#) should be called until it returns a null pointer, to allow libpq to process the error information completely.

Using [PQsendQuery](#) and [PQgetResult](#) solves one of [PQexec](#)'s problems: If a command string contains multiple SQL commands, the results of those commands can be obtained individually. (This allows a simple form of overlapped processing, by the way: the client can be handling the results of one command while the server is still working on later queries in the same command string.)

Another frequently-desired feature that can be obtained with [PQsendQuery](#) and [PQgetResult](#) is retrieving large query results a limited number of rows at a time. This is discussed in [Section 35.6](#).

By itself, calling [PQgetResult](#) will still cause the client to block until the server completes the next SQL command. This can be avoided by proper use of two more functions:

PQconsumeInput

If input is available from the server, consume it.


```
int PQconsumeInput(PGconn *conn);
```

`PQconsumeInput` normally returns 1 indicating “no error”, but returns 0 if there was some kind of trouble (in which case `PQerrorMessage` can be consulted). Note that the result does not say whether any input data was actually collected. After calling `PQconsumeInput`, the application can check `PQisBusy` and/or `PQnotifies` to see if their state has changed.

`PQconsumeInput` can be called even if the application is not prepared to deal with a result or notification just yet. The function will read available data and save it in a buffer, thereby causing a `select()` read-ready indication to go away. The application can thus use `PQconsumeInput` to clear the `select()` condition immediately, and then examine the results at leisure.

`PQisBusy`

Returns 1 if a command is busy, that is, `PQgetResult` would block waiting for input. A 0 return indicates that `PQgetResult` can be called with assurance of not blocking.

```
int PQisBusy(PGconn *conn);
```

`PQisBusy` will not itself attempt to read data from the server; therefore `PQconsumeInput` must be invoked first, or the busy state will never end.

A typical application using these functions will have a main loop that uses `select()` or `poll()` to wait for all the conditions that it must respond to. One of the conditions will be input available from the server, which in terms of `select()` means readable data on the file descriptor identified by `PQsocket`. When the main loop detects input ready, it should call `PQconsumeInput` to read the input. It can then call `PQisBusy`, followed by `PQgetResult` if `PQisBusy` returns false (0). It can also call `PQnotifies` to detect NOTIFY messages (see [Section 35.9](#)).

A client that uses `PQsendQuery/PQgetResult` can also attempt to cancel a command that is still being processed by the server; see [Section 35.7](#). But regardless of the return value of `PQcancelBlocking`, the application must continue with the normal result-reading sequence using `PQgetResult`. A successful cancellation will simply cause the command to terminate sooner than it would have otherwise.

By using the functions described above, it is possible to avoid blocking while waiting for input from the database server. However, it is still possible that the application will block waiting to send output to the server. This is relatively uncommon but can happen if very long SQL commands or data values are sent. (It is much more probable if the application sends data via `COPY IN`, however.) To prevent this possibility and achieve completely nonblocking database operation, the following additional functions can be used.

`PQsetnonblocking`

Sets the nonblocking status of the connection.

```
int PQsetnonblocking(PGconn *conn, int arg);
```

Sets the state of the connection to nonblocking if `arg` is 1, or blocking if `arg` is 0. Returns 0 if OK, -1 if error.

In the nonblocking state, successful calls to `PQsendQuery`, `PQputline`, `PQputnbytes`, `PQputCopyData`, and `PQendcopy` will not block; their changes are stored in the local output buffer until they are flushed. Unsuccessful calls will return an error and must be retried.

Note that `PQexec` does not honor nonblocking mode; if it is called, it will act in blocking fashion anyway.

`PQisnonblocking`

Returns the blocking status of the database connection.

```
int PQisnonblocking(const PGconn *conn);
```

Returns 1 if the connection is set to nonblocking mode and 0 if blocking.

PQflush

Attempts to flush any queued output data to the server. Returns 0 if successful (or if the send queue is empty), -1 if it failed for some reason, or 1 if it was unable to send all the data in the send queue yet (this case can only occur if the connection is nonblocking).

```
int PQflush(PGconn *conn);
```

After sending any command or data on a nonblocking connection, call `PQflush`. If it returns 1, wait for the socket to become read- or write-ready. If it becomes write-ready, call `PQflush` again. If it becomes read-ready, call `PQconsumeInput`, then call `PQflush` again. Repeat until `PQflush` returns 0. (It is necessary to check for read-ready and drain the input with `PQconsumeInput`, because the server can block trying to send us data, e.g., NOTICE messages, and won't read our data until we read its.) Once `PQflush` returns 0, wait for the socket to be read-ready and then read the response as described above.

35.5. Pipeline Mode

libpq pipeline mode allows applications to send a query without having to read the result of the previously sent query. Taking advantage of the pipeline mode, a client will wait less for the server, since multiple queries/results can be sent/received in a single network transaction.

While pipeline mode provides a significant performance boost, writing clients using the pipeline mode is more complex because it involves managing a queue of pending queries and finding which result corresponds to which query in the queue.

Pipeline mode also generally consumes more memory on both the client and server, though careful and aggressive management of the send/receive queue can mitigate this. This applies whether or not the connection is in blocking or non-blocking mode.

While libpq's pipeline API was introduced in Postgres Pro 14, it is a client-side feature which doesn't require special server support and works on any server that supports the v3 extended query protocol. For more information see [Section 56.2.4](#).

35.5.1. Using Pipeline Mode

To issue pipelines, the application must switch the connection into pipeline mode, which is done with `PQenterPipelineMode`. `PQpipelineStatus` can be used to test whether pipeline mode is active. In pipeline mode, only [asynchronous operations](#) that utilize the extended query protocol are permitted, command strings containing multiple SQL commands are disallowed, and so is `COPY`. Using synchronous command execution functions such as `PQfn`, `PQexec`, `PQexecParams`, `PQprepare`, `PQexecPrepared`, `PQdescribePrepared`, `PQdescribePortal`, `PQclosePrepared`, `PQclosePortal`, is an error condition. `PQsendQuery` is also disallowed, because it uses the simple query protocol. Once all dispatched commands have had their results processed, and the end pipeline result has been consumed, the application may return to non-pipelined mode with `PQexitPipelineMode`.

Note

It is best to use pipeline mode with libpq in [non-blocking mode](#). If used in blocking mode it is possible for a client/server deadlock to occur.¹

35.5.1.1. Issuing Queries

After entering pipeline mode, the application dispatches requests using `PQsendQueryParams` or its prepared-query sibling `PQsendQueryPrepared`. These requests are queued on the client-side until flushed to the server; this occurs when `PQpipelineSync` is used to establish a synchronization point in the pipeline, or when `PQflush` is called. The functions `PQsendPrepare`, `PQsendDescribePrepared`, `PQsend-`

¹ The client will block trying to send queries to the server, but the server will block trying to send results to the client from queries it has already processed. This only occurs when the client sends enough queries to fill both its output buffer and the server's receive buffer before it switches to processing input from the server, but it's hard to predict exactly when that will happen.

`DescribePortal`, `PQsendClosePrepared`, and `PQsendClosePortal` also work in pipeline mode. Result processing is described below.

The server executes statements, and returns results, in the order the client sends them. The server will begin executing the commands in the pipeline immediately, not waiting for the end of the pipeline. Note that results are buffered on the server side; the server flushes that buffer when a synchronization point is established with either `PQpipelineSync` or `PQsendPipelineSync`, or when `PQsendFlushRequest` is called. If any statement encounters an error, the server aborts the current transaction and does not execute any subsequent command in the queue until the next synchronization point; a `PGRES_PIPELINE_ABORTED` result is produced for each such command. (This remains true even if the commands in the pipeline would rollback the transaction.) Query processing resumes after the synchronization point.

It's fine for one operation to depend on the results of a prior one; for example, one query may define a table that the next query in the same pipeline uses. Similarly, an application may create a named prepared statement and execute it with later statements in the same pipeline.

35.5.1.2. Processing Results

To process the result of one query in a pipeline, the application calls `PQgetResult` repeatedly and handles each result until `PQgetResult` returns null. The result from the next query in the pipeline may then be retrieved using `PQgetResult` again and the cycle repeated. The application handles individual statement results as normal. When the results of all the queries in the pipeline have been returned, `PQgetResult` returns a result containing the status value `PGRES_PIPELINE_SYNC`.

The client may choose to defer result processing until the complete pipeline has been sent, or interleave that with sending further queries in the pipeline; see [Section 35.5.1.4](#).

`PQgetResult` behaves the same as for normal asynchronous processing except that it may contain the new `PGresult` types `PGRES_PIPELINE_SYNC` and `PGRES_PIPELINE_ABORTED`. `PGRES_PIPELINE_SYNC` is reported exactly once for each `PQpipelineSync` or `PQsendPipelineSync` at the corresponding point in the pipeline. `PGRES_PIPELINE_ABORTED` is emitted in place of a normal query result for the first error and all subsequent results until the next `PGRES_PIPELINE_SYNC`; see [Section 35.5.1.3](#).

`PQisBusy`, `PQconsumeInput`, etc operate as normal when processing pipeline results. In particular, a call to `PQisBusy` in the middle of a pipeline returns 0 if the results for all the queries issued so far have been consumed.

libpq does not provide any information to the application about the query currently being processed (except that `PQgetResult` returns null to indicate that we start returning the results of next query). The application must keep track of the order in which it sent queries, to associate them with their corresponding results. Applications will typically use a state machine or a FIFO queue for this.

35.5.1.3. Error Handling

From the client's perspective, after `PQresultStatus` returns `PGRES_FATAL_ERROR`, the pipeline is flagged as aborted. `PQresultStatus` will report a `PGRES_PIPELINE_ABORTED` result for each remaining queued operation in an aborted pipeline. The result for `PQpipelineSync` or `PQsendPipelineSync` is reported as `PGRES_PIPELINE_SYNC` to signal the end of the aborted pipeline and resumption of normal result processing.

The client *must* process results with `PQgetResult` during error recovery.

If the pipeline used an implicit transaction, then operations that have already executed are rolled back and operations that were queued to follow the failed operation are skipped entirely. The same behavior holds if the pipeline starts and commits a single explicit transaction (i.e. the first statement is `BEGIN` and the last is `COMMIT`) except that the session remains in an aborted transaction state at the end of the pipeline. If a pipeline contains *multiple explicit transactions*, all transactions that committed prior to the error remain committed, the currently in-progress transaction is aborted, and all subsequent operations are skipped completely, including subsequent transactions. If a pipeline synchronization point occurs

with an explicit transaction block in aborted state, the next pipeline will become aborted immediately unless the next command puts the transaction in normal mode with `ROLLBACK`.

Note

The client must not assume that work is committed when it *sends* a `COMMIT` — only when the corresponding result is received to confirm the commit is complete. Because errors arrive asynchronously, the application needs to be able to restart from the last *received* committed change and resend work done after that point if something goes wrong.

35.5.1.4. Interleaving Result Processing and Query Dispatch

To avoid deadlocks on large pipelines the client should be structured around a non-blocking event loop using operating system facilities such as `select`, `poll`, `WaitForMultipleObjectEx`, etc.

The client application should generally maintain a queue of work remaining to be dispatched and a queue of work that has been dispatched but not yet had its results processed. When the socket is writable it should dispatch more work. When the socket is readable it should read results and process them, matching them up to the next entry in its corresponding results queue. Based on available memory, results from the socket should be read frequently: there's no need to wait until the pipeline end to read the results. Pipelines should be scoped to logical units of work, usually (but not necessarily) one transaction per pipeline. There's no need to exit pipeline mode and re-enter it between pipelines, or to wait for one pipeline to finish before sending the next.

35.5.2. Functions Associated with Pipeline Mode

`PQpipelineStatus`

Returns the current pipeline mode status of the libpq connection.

```
PGpipelineStatus PQpipelineStatus(const PGconn *conn);
```

`PQpipelineStatus` can return one of the following values:

`PQ_PIPELINE_ON`

The libpq connection is in pipeline mode.

`PQ_PIPELINE_OFF`

The libpq connection is *not* in pipeline mode.

`PQ_PIPELINE_ABORTED`

The libpq connection is in pipeline mode and an error occurred while processing the current pipeline. The aborted flag is cleared when `PQgetResult` returns a result of type `PGRES_PIPELINE_SYNC`.

`PQenterPipelineMode`

Causes a connection to enter pipeline mode if it is currently idle or already in pipeline mode.

```
int PQenterPipelineMode(PGconn *conn);
```

Returns 1 for success. Returns 0 and has no effect if the connection is not currently idle, i.e., it has a result ready, or it is waiting for more input from the server, etc. This function does not actually send anything to the server, it just changes the libpq connection state.

`PQexitPipelineMode`

Causes a connection to exit pipeline mode if it is currently in pipeline mode with an empty queue and no pending results.

```
int PQexitPipelineMode(PGconn *conn);
```

Returns 1 for success. Returns 1 and takes no action if not in pipeline mode. If the current statement isn't finished processing, or `PQgetResult` has not been called to collect results from all previously sent query, returns 0 (in which case, use `PQerrorMessage` to get more information about the failure).

`PQpipelineSync`

Marks a synchronization point in a pipeline by sending a [sync message](#) and flushing the send buffer. This serves as the delimiter of an implicit transaction and an error recovery point; see [Section 35.5.1.3](#).

```
int PQpipelineSync(PGconn *conn);
```

Returns 1 for success. Returns 0 if the connection is not in pipeline mode or sending a [sync message](#) failed.

`PQsendPipelineSync`

Marks a synchronization point in a pipeline by sending a [sync message](#) without flushing the send buffer. This serves as the delimiter of an implicit transaction and an error recovery point; see [Section 35.5.1.3](#).

```
int PQsendPipelineSync(PGconn *conn);
```

Returns 1 for success. Returns 0 if the connection is not in pipeline mode or sending a [sync message](#) failed. Note that the message is not itself flushed to the server automatically; use `PQflush` if necessary.

`PQsendFlushRequest`

Sends a request for the server to flush its output buffer.

```
int PQsendFlushRequest(PGconn *conn);
```

Returns 1 for success. Returns 0 on any failure.

The server flushes its output buffer automatically as a result of `PQpipelineSync` being called, or on any request when not in pipeline mode; this function is useful to cause the server to flush its output buffer in pipeline mode without establishing a synchronization point. Note that the request is not itself flushed to the server automatically; use `PQflush` if necessary.

35.5.3. When to Use Pipeline Mode

Much like asynchronous query mode, there is no meaningful performance overhead when using pipeline mode. It increases client application complexity, and extra caution is required to prevent client/server deadlocks, but pipeline mode can offer considerable performance improvements, in exchange for increased memory usage from leaving state around longer.

Pipeline mode is most useful when the server is distant, i.e., network latency (“ping time”) is high, and also when many small operations are being performed in rapid succession. There is usually less benefit in using pipelined commands when each query takes many multiples of the client/server round-trip time to execute. A 100-statement operation run on a server 300 ms round-trip-time away would take 30 seconds in network latency alone without pipelining; with pipelining it may spend as little as 0.3 s waiting for results from the server.

Use pipelined commands when your application does lots of small `INSERT`, `UPDATE` and `DELETE` operations that can't easily be transformed into operations on sets, or into a `COPY` operation.

Pipeline mode is not useful when information from one operation is required by the client to produce the next operation. In such cases, the client would have to introduce a synchronization point and wait for a full client/server round-trip to get the results it needs. However, it's often possible to adjust the client design to exchange the required information server-side. Read-modify-write cycles are especially good candidates; for example:

```
BEGIN;  
SELECT x FROM mytable WHERE id = 42 FOR UPDATE;  
-- result: x=2  
-- client adds 1 to x:  
UPDATE mytable SET x = 3 WHERE id = 42;  
COMMIT;
```

could be much more efficiently done with:

```
UPDATE mytable SET x = x + 1 WHERE id = 42;
```

Pipelining is less useful, and more complex, when a single pipeline contains multiple transactions (see [Section 35.5.1.3](#)).

35.6. Retrieving Query Results in Chunks

Ordinarily, libpq collects an SQL command's entire result and returns it to the application as a single `PGresult`. This can be unworkable for commands that return a large number of rows. For such cases, applications can use `PQsendQuery` and `PQgetResult` in *single-row mode* or *chunked mode*. In these modes, result row(s) are returned to the application as they are received from the server, one at a time for single-row mode or in groups for chunked mode.

To enter one of these modes, call `PQsetSingleRowMode` or `PQsetChunkedRowsMode` immediately after a successful call of `PQsendQuery` (or a sibling function). This mode selection is effective only for the currently executing query. Then call `PQgetResult` repeatedly, until it returns null, as documented in [Section 35.4](#). If the query returns any rows, they are returned as one or more `PGresult` objects, which look like normal query results except for having status code `PGRES_SINGLE_TUPLE` for single-row mode or `PGRES_TUPLES_CHUNK` for chunked mode, instead of `PGRES_TUPLES_OK`. There is exactly one result row in each `PGRES_SINGLE_TUPLE` object, while a `PGRES_TUPLES_CHUNK` object contains at least one row but not more than the specified number of rows per chunk. After the last row, or immediately if the query returns zero rows, a zero-row object with status `PGRES_TUPLES_OK` is returned; this is the signal that no more rows will arrive. (But note that it is still necessary to continue calling `PQgetResult` until it returns null.) All of these `PGresult` objects will contain the same row description data (column names, types, etc.) that an ordinary `PGresult` object for the query would have. Each object should be freed with `PQclear` as usual.

When using pipeline mode, single-row or chunked mode needs to be activated for each query in the pipeline before retrieving results for that query with `PQgetResult`. See [Section 35.5](#) for more information.

`PQsetSingleRowMode`

Select single-row mode for the currently-executing query.

```
int PQsetSingleRowMode(PGconn *conn);
```

This function can only be called immediately after `PQsendQuery` or one of its sibling functions, before any other operation on the connection such as `PQconsumeInput` or `PQgetResult`. If called at the correct time, the function activates single-row mode for the current query and returns 1. Otherwise the mode stays unchanged and the function returns 0. In any case, the mode reverts to normal after completion of the current query.

`PQsetChunkedRowsMode`

Select chunked mode for the currently-executing query.

```
int PQsetChunkedRowsMode(PGconn *conn, int chunkSize);
```

This function is similar to `PQsetSingleRowMode`, except that it specifies retrieval of up to `chunkSize` rows per `PGresult`, not necessarily just one row. This function can only be called immediately after `PQsendQuery` or one of its sibling functions, before any other operation on the connection such as

`PQconsumeInput` or `PQgetResult`. If called at the correct time, the function activates chunked mode for the current query and returns 1. Otherwise the mode stays unchanged and the function returns 0. In any case, the mode reverts to normal after completion of the current query.

Caution

While processing a query, the server may return some rows and then encounter an error, causing the query to be aborted. Ordinarily, libpq discards any such rows and reports only the error. But in single-row or chunked mode, some rows may have already been returned to the application. Hence, the application will see some `PGRES_SINGLE_TUPLE` or `PGRES_TUPLES_CHUNK` `PGresult` objects followed by a `PGRES_FATAL_ERROR` object. For proper transactional behavior, the application must be designed to discard or undo whatever has been done with the previously-processed rows, if the query ultimately fails.

35.7. Canceling Queries in Progress

35.7.1. Functions for Sending Cancel Requests

`PQcancelCreate`

Prepares a connection over which a cancel request can be sent.

```
PGcancelConn *PQcancelCreate(PGconn *conn);
```

`PQcancelCreate` creates a `PGcancelConn` object, but it won't instantly start sending a cancel request over this connection. A cancel request can be sent over this connection in a blocking manner using `PQcancelBlocking` and in a non-blocking manner using `PQcancelStart`. The return value can be passed to `PQcancelStatus` to check if the `PGcancelConn` object was created successfully. The `PGcancelConn` object is an opaque structure that is not meant to be accessed directly by the application. This `PGcancelConn` object can be used to cancel the query that's running on the original connection in a thread-safe way.

Many connection parameters of the original client will be reused when setting up the connection for the cancel request. Importantly, if the original connection requires encryption of the connection and/or verification of the target host (using `sslmode` or `gssencmode`), then the connection for the cancel request is made with these same requirements. Any connection options that are only used during authentication or after authentication of the client are ignored though, because cancellation requests do not require authentication and the connection is closed right after the cancellation request is submitted.

Note that when `PQcancelCreate` returns a non-null pointer, you must call `PQcancelFinish` when you are finished with it, in order to dispose of the structure and any associated memory blocks. This must be done even if the cancel request failed or was abandoned.

`PQcancelBlocking`

Requests that the server abandons processing of the current command in a blocking manner.

```
int PQcancelBlocking(PGcancelConn *cancelConn);
```

The request is made over the given `PGcancelConn`, which needs to be created with `PQcancelCreate`. The return value of `PQcancelBlocking` is 1 if the cancel request was successfully dispatched and 0 if not. If it was unsuccessful, the error message can be retrieved using `PQcancelErrorMessage`.

Successful dispatch of the cancellation is no guarantee that the request will have any effect, however. If the cancellation is effective, the command being canceled will terminate early and return an error result. If the cancellation fails (say, because the server was already done processing the command), then there will be no visible result at all.

PQcancelStart
PQcancelPoll

Requests that the server abandons processing of the current command in a non-blocking manner.

```
int PQcancelStart(PGcancelConn *cancelConn);
```

```
PostgresPollingStatusType PQcancelPoll(PGcancelConn *cancelConn);
```

The request is made over the given `PGcancelConn`, which needs to be created with `PQcancelCreate`. The return value of `PQcancelStart` is 1 if the cancellation request could be started and 0 if not. If it was unsuccessful, the error message can be retrieved using `PQcancelErrorMessage`.

If `PQcancelStart` succeeds, the next stage is to poll libpq so that it can proceed with the cancel connection sequence. Use `PQcancelSocket` to obtain the descriptor of the socket underlying the database connection. (Caution: do not assume that the socket remains the same across `PQcancelPoll` calls.) Loop thus: If `PQcancelPoll(cancelConn)` last returned `PGRES_POLLING_READING`, wait until the socket is ready to read (as indicated by `select()`, `poll()`, or similar system function). Then call `PQcancelPoll(cancelConn)` again. Conversely, if `PQcancelPoll(cancelConn)` last returned `PGRES_POLLING_WRITING`, wait until the socket is ready to write, then call `PQcancelPoll(cancelConn)` again. On the first iteration, i.e., if you have yet to call `PQcancelPoll(cancelConn)`, behave as if it last returned `PGRES_POLLING_WRITING`. Continue this loop until `PQcancelPoll(cancelConn)` returns `PGRES_POLLING_FAILED`, indicating the connection procedure has failed, or `PGRES_POLLING_OK`, indicating cancel request was successfully dispatched.

Successful dispatch of the cancellation is no guarantee that the request will have any effect, however. If the cancellation is effective, the command being canceled will terminate early and return an error result. If the cancellation fails (say, because the server was already done processing the command), then there will be no visible result at all.

At any time during connection, the status of the connection can be checked by calling `PQcancelStatus`. If this call returns `CONNECTION_BAD`, then the cancel procedure has failed; if the call returns `CONNECTION_OK`, then cancel request was successfully dispatched. Both of these states are equally detectable from the return value of `PQcancelPoll`, described above. Other states might also occur during (and only during) an asynchronous connection procedure. These indicate the current stage of the connection procedure and might be useful to provide feedback to the user for example. These statuses are:

`CONNECTION_ALLOCATED`

Waiting for a call to `PQcancelStart` or `PQcancelBlocking`, to actually open the socket. This is the connection state right after calling `PQcancelCreate` or `PQcancelReset`. No connection to the server has been initiated yet at this point. To actually start sending the cancel request use `PQcancelStart` or `PQcancelBlocking`.

`CONNECTION_STARTED`

Waiting for connection to be made.

`CONNECTION_MADE`

Connection OK; waiting to send.

`CONNECTION_AWAITING_RESPONSE`

Waiting for a response from the server.

`CONNECTION_SSL_STARTUP`

Negotiating SSL encryption.

`CONNECTION_GSS_STARTUP`

Negotiating GSS encryption.

Note that, although these constants will remain (in order to maintain compatibility), an application should never rely upon these occurring in a particular order, or at all, or on the status always being one of these documented values. An application might do something like this:

```
switch (PQcancelStatus(conn))
{
    case CONNECTION_STARTED:
        feedback = "Connecting...";
        break;

    case CONNECTION_MADE:
        feedback = "Connected to server...";
        break;

    .
    .
    .

    default:
        feedback = "Connecting...";
}
```

The `connect_timeout` connection parameter is ignored when using `PQcancelPoll`; it is the application's responsibility to decide whether an excessive amount of time has elapsed. Otherwise, `PQcancelStart` followed by a `PQcancelPoll` loop is equivalent to [PQcancelBlocking](#).

`PQcancelStatus`

Returns the status of the cancel connection.

```
ConnStatusType PQcancelStatus(const PGcancelConn *cancelConn);
```

The status can be one of a number of values. However, only three of these are seen outside of an asynchronous cancel procedure: `CONNECTION_ALLOCATED`, `CONNECTION_OK` and `CONNECTION_BAD`. The initial state of a `PGcancelConn` that's successfully created using [PQcancelCreate](#) is `CONNECTION_ALLOCATED`. A cancel request that was successfully dispatched has the status `CONNECTION_OK`. A failed cancel attempt is signaled by status `CONNECTION_BAD`. An OK status will remain so until [PQcancelFinish](#) or [PQcancelReset](#) is called.

See the entry for [PQcancelStart](#) with regards to other status codes that might be returned.

Successful dispatch of the cancellation is no guarantee that the request will have any effect, however. If the cancellation is effective, the command being canceled will terminate early and return an error result. If the cancellation fails (say, because the server was already done processing the command), then there will be no visible result at all.

`PQcancelSocket`

Obtains the file descriptor number of the cancel connection socket to the server.

```
int PQcancelSocket(const PGcancelConn *cancelConn);
```

A valid descriptor will be greater than or equal to 0; a result of -1 indicates that no server connection is currently open. This might change as a result of calling any of the functions in this section on the `PGcancelConn` (except for [PQcancelErrorMessage](#) and `PQcancelSocket` itself).

`PQcancelErrorMessage`

Returns the error message most recently generated by an operation on the cancel connection.

```
char *PQcancelErrorMessage(const PGcancelConn *cancelconn);
```

Nearly all libpq functions that take a `PGcancelConn` will set a message for [PQcancelErrorMessage](#) if they fail. Note that by libpq convention, a nonempty [PQcancelErrorMessage](#) result can consist

of multiple lines, and will include a trailing newline. The caller should not free the result directly. It will be freed when the associated `PGcancelConn` handle is passed to `PQcancelFinish`. The result string should not be expected to remain the same across operations on the `PGcancelConn` structure.

`PQcancelFinish`

Closes the cancel connection (if it did not finish sending the cancel request yet). Also frees memory used by the `PGcancelConn` object.

```
void PQcancelFinish(PGcancelConn *cancelConn);
```

Note that even if the cancel attempt fails (as indicated by `PQcancelStatus`), the application should call `PQcancelFinish` to free the memory used by the `PGcancelConn` object. The `PGcancelConn` pointer must not be used again after `PQcancelFinish` has been called.

`PQcancelReset`

Resets the `PGcancelConn` so it can be reused for a new cancel connection.

```
void PQcancelReset(PGcancelConn *cancelConn);
```

If the `PGcancelConn` is currently used to send a cancel request, then this connection is closed. It will then prepare the `PGcancelConn` object such that it can be used to send a new cancel request.

This can be used to create one `PGcancelConn` for a `PGconn` and reuse it multiple times throughout the lifetime of the original `PGconn`.

35.7.2. Obsolete Functions for Sending Cancel Requests

These functions represent older methods of sending cancel requests. Although they still work, they are deprecated due to not sending the cancel requests in an encrypted manner, even when the original connection specified `sslmode` or `gssencmode` to require encryption. Thus these older methods are heavily discouraged from being used in new code, and it is recommended to change existing code to use the new functions instead.

`PQgetCancel`

Creates a data structure containing the information needed to cancel a command using `PQcancel`.

```
PGcancel *PQgetCancel(PGconn *conn);
```

`PQgetCancel` creates a `PGcancel` object given a `PGconn` connection object. It will return `NULL` if the given `conn` is `NULL` or an invalid connection. The `PGcancel` object is an opaque structure that is not meant to be accessed directly by the application; it can only be passed to `PQcancel` or `PQfreeCancel`.

`PQfreeCancel`

Frees a data structure created by `PQgetCancel`.

```
void PQfreeCancel(PGcancel *cancel);
```

`PQfreeCancel` frees a data object previously created by `PQgetCancel`.

`PQcancel`

`PQcancel` is a deprecated and insecure variant of `PQcancelBlocking`, but one that can be used safely from within a signal handler.

```
int PQcancel(PGcancel *cancel, char *errbuf, int errbufsize);
```

`PQcancel` only exists because of backwards compatibility reasons. `PQcancelBlocking` should be used instead. The only benefit that `PQcancel` has is that it can be safely invoked from a signal handler, if the `errbuf` is a local variable in the signal handler. However, this is generally not considered a big enough benefit to be worth the security issues that this function has.

The `PGcancel` object is read-only as far as `PQcancel` is concerned, so it can also be invoked from a thread that is separate from the one manipulating the `PGconn` object.

The return value of `PQcancel` is 1 if the cancel request was successfully dispatched and 0 if not. If not, `errbuf` is filled with an explanatory error message. `errbuf` must be a char array of size `errbufsize` (the recommended size is 256 bytes).

`PQrequestCancel`

`PQrequestCancel` is a deprecated and insecure variant of `PQcancelBlocking`.

```
int PQrequestCancel(PGconn *conn);
```

`PQrequestCancel` only exists because of backwards compatibility reasons. `PQcancelBlocking` should be used instead. There is no benefit to using `PQrequestCancel` over `PQcancelBlocking`.

Requests that the server abandon processing of the current command. It operates directly on the `PGconn` object, and in case of failure stores the error message in the `PGconn` object (whence it can be retrieved by `PQerrorMessage`). Although the functionality is the same, this approach is not safe within multiple-thread programs or signal handlers, since it is possible that overwriting the `PGconn`'s error message will mess up the operation currently in progress on the connection.

35.8. The Fast-Path Interface

Postgres Pro provides a fast-path interface to send simple function calls to the server.

Tip

This interface is somewhat obsolete, as one can achieve similar performance and greater functionality by setting up a prepared statement to define the function call. Then, executing the statement with binary transmission of parameters and results substitutes for a fast-path function call.

The function `PQfnrequests` execution of a server function via the fast-path interface:

```
PGresult *PQfn(PGconn *conn,
               int fnid,
               int *result_buf,
               int *result_len,
               int result_is_int,
               const PQArgBlock *args,
               int nargs);
```

```
typedef struct
{
    int len;
    int isint;
    union
    {
        int *ptr;
        int integer;
    } u;
} PQArgBlock;
```

The `fnid` argument is the OID of the function to be executed. `args` and `nargs` define the parameters to be passed to the function; they must match the declared function argument list. When the `isint` field of a parameter structure is true, the `u.integer` value is sent to the server as an integer of the indicated length (this must be 2 or 4 bytes); proper byte-swapping occurs. When `isint` is false, the indicated number of bytes at `*u.ptr` are sent with no processing; the data must be in the format expected by the

server for binary transmission of the function's argument data type. (The declaration of `u.ptr` as being of type `int *` is historical; it would be better to consider it `void *`.) `result_buf` points to the buffer in which to place the function's return value. The caller must have allocated sufficient space to store the return value. (There is no check!) The actual result length in bytes will be returned in the integer pointed to by `result_len`. If a 2- or 4-byte integer result is expected, set `result_is_int` to 1, otherwise set it to 0. Setting `result_is_int` to 1 causes libpq to byte-swap the value if necessary, so that it is delivered as a proper `int` value for the client machine; note that a 4-byte integer is delivered into `*result_buf` for either allowed result size. When `result_is_int` is 0, the binary-format byte string sent by the server is returned unmodified. (In this case it's better to consider `result_buf` as being of type `void *`.)

PQfn always returns a valid PGresult pointer, with status PGRES_COMMAND_OK for success or PGRES_FATAL_ERROR if some problem was encountered. The result status should be checked before the result is used. The caller is responsible for freeing the PGresult with PQclear when it is no longer needed.

To pass a NULL argument to the function, set the `len` field of that parameter structure to -1; the `isint` and `u` fields are then irrelevant.

If the function returns NULL, `*result_len` is set to -1, and `*result_buf` is not modified.

Note that it is not possible to handle set-valued results when using this interface. Also, the function must be a plain function, not an aggregate, window function, or procedure.

35.9. Asynchronous Notification

Postgres Pro offers asynchronous notification via the LISTEN and NOTIFY commands. A client session registers its interest in a particular notification channel with the LISTEN command (and can stop listening with the UNLISTEN command). All sessions listening on a particular channel will be notified asynchronously when a NOTIFY command with that channel name is executed by any session. A “payload” string can be passed to communicate additional data to the listeners.

libpq applications submit LISTEN, UNLISTEN, and NOTIFY commands as ordinary SQL commands. The arrival of NOTIFY messages can subsequently be detected by calling PQnotifies.

The function PQnotifies returns the next notification from a list of unhandled notification messages received from the server. It returns a null pointer if there are no pending notifications. Once a notification is returned from PQnotifies, it is considered handled and will be removed from the list of notifications.

```
PGnotify *PQnotifies(PGconn *conn);
```

```
typedef struct pgNotify
{
    char *relname;           /* notification channel name */
    int  be_pid;             /* process ID of notifying server process */
    char *extra;             /* notification payload string */
} PGnotify;
```

After processing a PGnotify object returned by PQnotifies, be sure to free it with PQfreemem. It is sufficient to free the PGnotify pointer; the `relname` and `extra` fields do not represent separate allocations. (The names of these fields are historical; in particular, channel names need not have anything to do with relation names.)

Example 35.2 gives a sample program that illustrates the use of asynchronous notification.

PQnotifies does not actually read data from the server; it just returns messages previously absorbed by another libpq function. In ancient releases of libpq, the only way to ensure timely receipt of NOTIFY messages was to constantly submit commands, even empty ones, and then check PQnotifies after each PQexec. While this still works, it is deprecated as a waste of processing power.

A better way to check for NOTIFY messages when you have no useful commands to execute is to call PQconsumeInput, then check PQnotifies. You can use select() to wait for data to arrive from the server,

thereby using no CPU power unless there is something to do. (See [PQsocket](#) to obtain the file descriptor number to use with `select()`.) Note that this will work OK whether you submit commands with [PQsendQuery/PQgetResult](#) or simply use [PQexec](#). You should, however, remember to check `PQnotifies` after each [PQgetResult](#) or [PQexec](#), to see if any notifications came in during the processing of the command.

35.10. Functions Associated with the COPY Command

The `COPY` command in Postgres Pro has options to read from or write to the network connection used by libpq. The functions described in this section allow applications to take advantage of this capability by supplying or consuming copied data.

The overall process is that the application first issues the SQL `COPY` command via [PQexec](#) or one of the equivalent functions. The response to this (if there is no error in the command) will be a `PGresult` object bearing a status code of `PGRES_COPY_OUT` or `PGRES_COPY_IN` (depending on the specified copy direction). The application should then use the functions of this section to receive or transmit data rows. When the data transfer is complete, another `PGresult` object is returned to indicate success or failure of the transfer. Its status will be `PGRES_COMMAND_OK` for success or `PGRES_FATAL_ERROR` if some problem was encountered. At this point further SQL commands can be issued via [PQexec](#). (It is not possible to execute other SQL commands using the same connection while the `COPY` operation is in progress.)

If a `COPY` command is issued via [PQexec](#) in a string that could contain additional commands, the application must continue fetching results via [PQgetResult](#) after completing the `COPY` sequence. Only when [PQgetResult](#) returns `NULL` is it certain that the [PQexec](#) command string is done and it is safe to issue more commands.

The functions of this section should be executed only after obtaining a result status of `PGRES_COPY_OUT` or `PGRES_COPY_IN` from [PQexec](#) or [PQgetResult](#).

A `PGresult` object bearing one of these status values carries some additional data about the `COPY` operation that is starting. This additional data is available using functions that are also used in connection with query results:

`PQnfields`

Returns the number of columns (fields) to be copied.

`PQbinaryTuples`

0 indicates the overall copy format is textual (rows separated by newlines, columns separated by separator characters, etc.). 1 indicates the overall copy format is binary. See [COPY](#) for more information.

`PQfformat`

Returns the format code (0 for text, 1 for binary) associated with each column of the copy operation. The per-column format codes will always be zero when the overall copy format is textual, but the binary format can support both text and binary columns. (However, as of the current implementation of `COPY`, only binary columns appear in a binary copy; so the per-column formats always match the overall format at present.)

35.10.1. Functions for Sending COPY Data

These functions are used to send data during `COPY FROM STDIN`. They will fail if called when the connection is not in `COPY_IN` state.

`PQputCopyData`

Sends data to the server during `COPY_IN` state.

```
int PQputCopyData(PGconn *conn,
                  const char *buffer,
                  int nbytes);
```

Transmits the `COPY` data in the specified *buffer*, of length *nbytes*, to the server. The result is 1 if the data was queued, zero if it was not queued because of full buffers (this will only happen in nonblocking mode), or -1 if an error occurred. (Use [PQerrorMessage](#) to retrieve details if the return value is -1. If the value is zero, wait for write-ready and try again.)

The application can divide the `COPY` data stream into buffer loads of any convenient size. Buffer-load boundaries have no semantic significance when sending. The contents of the data stream must match the data format expected by the `COPY` command; see [COPY](#) for details.

PQputCopyEnd

Sends end-of-data indication to the server during `COPY_IN` state.

```
int PQputCopyEnd(PGconn *conn,
                 const char *errmsg);
```

Ends the `COPY_IN` operation successfully if *errmsg* is NULL. If *errmsg* is not NULL then the `COPY` is forced to fail, with the string pointed to by *errmsg* used as the error message. (One should not assume that this exact error message will come back from the server, however, as the server might have already failed the `COPY` for its own reasons.)

The result is 1 if the termination message was sent; or in nonblocking mode, this may only indicate that the termination message was successfully queued. (In nonblocking mode, to be certain that the data has been sent, you should next wait for write-ready and call [PQflush](#), repeating until it returns zero.) Zero indicates that the function could not queue the termination message because of full buffers; this will only happen in nonblocking mode. (In this case, wait for write-ready and try the [PQputCopyEnd](#) call again.) If a hard error occurs, -1 is returned; you can use [PQerrorMessage](#) to retrieve details.

After successfully calling [PQputCopyEnd](#), call [PQgetResult](#) to obtain the final result status of the `COPY` command. One can wait for this result to be available in the usual way. Then return to normal operation.

35.10.2. Functions for Receiving `COPY` Data

These functions are used to receive data during `COPY TO STDOUT`. They will fail if called when the connection is not in `COPY_OUT` state.

PQgetCopyData

Receives data from the server during `COPY_OUT` state.

```
int PQgetCopyData(PGconn *conn,
                 char **buffer,
                 int async);
```

Attempts to obtain another row of data from the server during a `COPY`. Data is always returned one data row at a time; if only a partial row is available, it is not returned. Successful return of a data row involves allocating a chunk of memory to hold the data. The *buffer* parameter must be non-NULL. **buffer* is set to point to the allocated memory, or to NULL in cases where no buffer is returned. A non-NULL result buffer should be freed using [PQfreemem](#) when no longer needed.

When a row is successfully returned, the return value is the number of data bytes in the row (this will always be greater than zero). The returned string is always null-terminated, though this is probably only useful for textual `COPY`. A result of zero indicates that the `COPY` is still in progress, but no row is yet available (this is only possible when *async* is true). A result of -1 indicates that the `COPY` is done. A result of -2 indicates that an error occurred (consult [PQerrorMessage](#) for the reason).

When *async* is true (not zero), [PQgetCopyData](#) will not block waiting for input; it will return zero if the `COPY` is still in progress but no complete row is available. (In this case wait for read-ready

and then call `PQconsumeInput` before calling `PQgetCopyData` again.) When `async` is false (zero), `PQgetCopyData` will block until data is available or the operation completes.

After `PQgetCopyData` returns -1, call `PQgetResult` to obtain the final result status of the `COPY` command. One can wait for this result to be available in the usual way. Then return to normal operation.

35.10.3. Obsolete Functions for COPY

These functions represent older methods of handling `COPY`. Although they still work, they are deprecated due to poor error handling, inconvenient methods of detecting end-of-data, and lack of support for binary or nonblocking transfers.

`PQgetline`

Reads a newline-terminated line of characters (transmitted by the server) into a buffer string of size `length`.

```
int PQgetline(PGconn *conn,
              char *buffer,
              int length);
```

This function copies up to `length-1` characters into the buffer and converts the terminating newline into a zero byte. `PQgetline` returns `EOF` at the end of input, 0 if the entire line has been read, and 1 if the buffer is full but the terminating newline has not yet been read.

Note that the application must check to see if a new line consists of the two characters `\.`, which indicates that the server has finished sending the results of the `COPY` command. If the application might receive lines that are more than `length-1` characters long, care is needed to be sure it recognizes the `\.` line correctly (and does not, for example, mistake the end of a long data line for a terminator line).

`PQgetlineAsync`

Reads a row of `COPY` data (transmitted by the server) into a buffer without blocking.

```
int PQgetlineAsync(PGconn *conn,
                  char *buffer,
                  int bufsize);
```

This function is similar to `PQgetline`, but it can be used by applications that must read `COPY` data asynchronously, that is, without blocking. Having issued the `COPY` command and gotten a `PGRES_COPY_OUT` response, the application should call `PQconsumeInput` and `PQgetlineAsync` until the end-of-data signal is detected.

Unlike `PQgetline`, this function takes responsibility for detecting end-of-data.

On each call, `PQgetlineAsync` will return data if a complete data row is available in libpq's input buffer. Otherwise, no data is returned until the rest of the row arrives. The function returns -1 if the end-of-copy-data marker has been recognized, or 0 if no data is available, or a positive number giving the number of bytes of data returned. If -1 is returned, the caller must next call `PQendcopy`, and then return to normal processing.

The data returned will not extend beyond a data-row boundary. If possible a whole row will be returned at one time. But if the buffer offered by the caller is too small to hold a row sent by the server, then a partial data row will be returned. With textual data this can be detected by testing whether the last returned byte is `\n` or not. (In a binary `COPY`, actual parsing of the `COPY` data format will be needed to make the equivalent determination.) The returned string is not null-terminated. (If you want to add a terminating null, be sure to pass a `bufsize` one smaller than the room actually available.)

`PQputline`

Sends a null-terminated string to the server. Returns 0 if OK and `EOF` if unable to send the string.

```
int PQputline(PGconn *conn,
              const char *string);
```

The `COPY` data stream sent by a series of calls to `PQputline` has the same format as that returned by `PQgetlineAsync`, except that applications are not obliged to send exactly one data row per `PQputline` call; it is okay to send a partial line or multiple lines per call.

Note

Before Postgres Pro protocol 3.0, it was necessary for the application to explicitly send the two characters `\.` as a final line to indicate to the server that it had finished sending `COPY` data. While this still works, it is deprecated and the special meaning of `\.` can be expected to be removed in a future release. It is sufficient to call `PQendcopy` after having sent the actual data.

PQputnbytes

Sends a non-null-terminated string to the server. Returns 0 if OK and `EOF` if unable to send the string.

```
int PQputnbytes(PGconn *conn,
                const char *buffer,
                int nbytes);
```

This is exactly like `PQputline`, except that the data buffer need not be null-terminated since the number of bytes to send is specified directly. Use this procedure when sending binary data.

PQendcopy

Synchronizes with the server.

```
int PQendcopy(PGconn *conn);
```

This function waits until the server has finished the copying. It should either be issued when the last string has been sent to the server using `PQputline` or when the last string has been received from the server using `PQgetline`. It must be issued or the server will get “out of sync” with the client. Upon return from this function, the server is ready to receive the next SQL command. The return value is 0 on successful completion, nonzero otherwise. (Use `PQerrorMessage` to retrieve details if the return value is nonzero.)

When using `PQgetResult`, the application should respond to a `PGRES_COPY_OUT` result by executing `PQgetline` repeatedly, followed by `PQendcopy` after the terminator line is seen. It should then return to the `PQgetResult` loop until `PQgetResult` returns a null pointer. Similarly a `PGRES_COPY_IN` result is processed by a series of `PQputline` calls followed by `PQendcopy`, then return to the `PQgetResult` loop. This arrangement will ensure that a `COPY` command embedded in a series of SQL commands will be executed correctly.

Older applications are likely to submit a `COPY` via `PQexec` and assume that the transaction is done after `PQendcopy`. This will work correctly only if the `COPY` is the only SQL command in the command string.

35.11. Control Functions

These functions control miscellaneous details of libpq's behavior.

PQclientEncoding

Returns the client encoding.

```
int PQclientEncoding(const PGconn *conn);
```

Note that it returns the encoding ID, not a symbolic string such as `EUC_JP`. If unsuccessful, it returns -1. To convert an encoding ID to an encoding name, you can use:

```
char *pg_encoding_to_char(int encoding_id);
```

PQsetClientEncoding

Sets the client encoding.

```
int PQsetClientEncoding(PGconn *conn, const char *encoding);
```

conn is a connection to the server, and *encoding* is the encoding you want to use. If the function successfully sets the encoding, it returns 0, otherwise -1. The current encoding for this connection can be determined by using [PQclientEncoding](#).

PQsetErrorVerbosity

Determines the verbosity of messages returned by [PQerrorMessage](#) and [PQresultErrorMessage](#).

```
typedef enum
```

```
{
    PQERRORS_TERSE,
    PQERRORS_DEFAULT,
    PQERRORS_VERBOSE,
    PQERRORS_SQLSTATE
} PGVerbosity;
```

```
PGVerbosity PQsetErrorVerbosity(PGconn *conn, PGVerbosity verbosity);
```

[PQsetErrorVerbosity](#) sets the verbosity mode, returning the connection's previous setting. In *TERSE* mode, returned messages include severity, primary text, and position only; this will normally fit on a single line. The *DEFAULT* mode produces messages that include the above plus any detail, hint, or context fields (these might span multiple lines). The *VERBOSE* mode includes all available fields. The *SQLSTATE* mode includes only the error severity and the *SQLSTATE* error code, if one is available (if not, the output is like *TERSE* mode).

Changing the verbosity setting does not affect the messages available from already-existing *PGresult* objects, only subsequently-created ones. (But see [PQresultVerboseErrorMessage](#) if you want to print a previous error with a different verbosity.)

PQsetErrorContextVisibility

Determines the handling of *CONTEXT* fields in messages returned by [PQerrorMessage](#) and [PQresultErrorMessage](#).

```
typedef enum
```

```
{
    PQSHOW_CONTEXT_NEVER,
    PQSHOW_CONTEXT_ERRORS,
    PQSHOW_CONTEXT_ALWAYS
} PGContextVisibility;
```

```
PGContextVisibility PQsetErrorContextVisibility(PGconn *conn, PGContextVisibility
show_context);
```

[PQsetErrorContextVisibility](#) sets the context display mode, returning the connection's previous setting. This mode controls whether the *CONTEXT* field is included in messages. The *NEVER* mode never includes *CONTEXT*, while *ALWAYS* always includes it if available. In *ERRORS* mode (the default), *CONTEXT* fields are included only in error messages, not in notices and warnings. (However, if the verbosity setting is *TERSE* or *SQLSTATE*, *CONTEXT* fields are omitted regardless of the context display mode.)

Changing this mode does not affect the messages available from already-existing *PGresult* objects, only subsequently-created ones. (But see [PQresultVerboseErrorMessage](#) if you want to print a previous error with a different display mode.)

PQtrace

Enables tracing of the client/server communication to a debugging file stream.


```
void PQtrace(PGconn *conn, FILE *stream);
```

Each line consists of: an optional timestamp, a direction indicator (**F** for messages from client to server or **B** for messages from server to client), message length, message type, and message contents. Non-message contents fields (timestamp, direction, length and message type) are separated by a tab. Message contents are separated by a space. Protocol strings are enclosed in double quotes, while strings used as data values are enclosed in single quotes. Non-printable chars are printed as hexadecimal escapes. Further message-type-specific detail can be found in [Section 56.7](#).

Note

On Windows, if the libpq library and an application are compiled with different flags, this function call will crash the application because the internal representation of the `FILE` pointers differ. Specifically, multithreaded/single-threaded, release/debug, and static/dynamic flags should be the same for the library and all applications using that library.

`PQsetTraceFlags`

Controls the tracing behavior of client/server communication.

```
void PQsetTraceFlags(PGconn *conn, int flags);
```

`flags` contains flag bits describing the operating mode of tracing. If `flags` contains `PQTRACE_SUPPRESS_TIMESTAMPS`, then the timestamp is not included when printing each message. If `flags` contains `PQTRACE_REGRESS_MODE`, then some fields are redacted when printing each message, such as object OIDs, to make the output more convenient to use in testing frameworks. This function must be called after calling `PQtrace`.

`PQuntrace`

Disables tracing started by `PQtrace`.

```
void PQuntrace(PGconn *conn);
```

35.12. Miscellaneous Functions

As always, there are some functions that just don't fit anywhere.

`PQfreemem`

Frees memory allocated by libpq.

```
void PQfreemem(void *ptr);
```

Frees memory allocated by libpq, particularly `PQescapeByteaConn`, `PQescapeBytea`, `PQunescapeBytea`, and `PQnotifies`. It is particularly important that this function, rather than `free()`, be used on Microsoft Windows. This is because allocating memory in a DLL and releasing it in the application works only if multithreaded/single-threaded, release/debug, and static/dynamic flags are the same for the DLL and the application. On non-Microsoft Windows platforms, this function is the same as the standard library function `free()`.

`PQconninfoFree`

Frees the data structures allocated by `PQconninfoDefaults` or `PQconninfoParse`.

```
void PQconninfoFree(PQconninfoOption *connOptions);
```

If the argument is a `NULL` pointer, no operation is performed.

A simple `PQfreemem` will not do for this, since the array contains references to subsidiary strings.

`PQencryptPasswordConn`

Prepares the encrypted form of a Postgres Pro password.


```
char *PQencryptPasswordConn(PGconn *conn, const char *passwd, const char *user,
    const char *algorithm);
```

This function is intended to be used by client applications that wish to send commands like `ALTER USER joe PASSWORD 'pwd'`. It is good practice not to send the original cleartext password in such a command, because it might be exposed in command logs, activity displays, and so on. Instead, use this function to convert the password to encrypted form before it is sent.

The *passwd* and *user* arguments are the cleartext password, and the SQL name of the user it is for. *algorithm* specifies the encryption algorithm to use to encrypt the password. Currently supported algorithms are `md5` and `scram-sha-256` (on and off are also accepted as aliases for `md5`, for compatibility with older server versions). Note that support for `scram-sha-256` was introduced in Postgres Pro version 10, and will not work correctly with older server versions. If *algorithm* is `NULL`, this function will query the server for the current value of the `password_encryption` setting. That can block, and will fail if the current transaction is aborted, or if the connection is busy executing another query. If you wish to use the default algorithm for the server but want to avoid blocking, query `password_encryption` yourself before calling `PQencryptPasswordConn`, and pass that value as the *algorithm*.

The return value is a string allocated by `malloc`. The caller can assume the string doesn't contain any special characters that would require escaping. Use `PQfreemem` to free the result when done with it. On error, returns `NULL`, and a suitable message is stored in the connection object.

`PQchangePassword`

Changes a Postgres Pro password.

```
PGresult *PQchangePassword(PGconn *conn, const char *user, const char *passwd);
```

This function uses `PQencryptPasswordConn` to build and execute the command `ALTER USER ... PASSWORD '...'`, thereby changing the user's password. It exists for the same reason as `PQencryptPasswordConn`, but is more convenient as it both builds and runs the command for you. `PQencryptPasswordConn` is passed a `NULL` for the algorithm argument, hence encryption is done according to the server's `password_encryption` setting.

The *user* and *passwd* arguments are the SQL name of the target user, and the new cleartext password.

Returns a `PGresult` pointer representing the result of the `ALTER USER` command, or a null pointer if the routine failed before issuing any command. The `PQresultStatus` function should be called to check the return value for any errors (including the value of a null pointer, in which case it will return `PGRES_FATAL_ERROR`). Use `PQerrorMessage` to get more information about such errors.

`PQencryptPassword`

Prepares the md5-encrypted form of a Postgres Pro password.

```
char *PQencryptPassword(const char *passwd, const char *user);
```

`PQencryptPassword` is an older, deprecated version of `PQencryptPasswordConn`. The difference is that `PQencryptPassword` does not require a connection object, and `md5` is always used as the encryption algorithm.

`PQmakeEmptyPGresult`

Constructs an empty `PGresult` object with the given status.

```
PGresult *PQmakeEmptyPGresult(PGconn *conn, ExecStatusType status);
```

This is libpq's internal function to allocate and initialize an empty `PGresult` object. This function returns `NULL` if memory could not be allocated. It is exported because some applications find it useful to generate result objects (particularly objects with error status) themselves. If *conn* is not null and *status* indicates an error, the current error message of the specified connection is copied into the `PGresult`. Also, if *conn* is not null, any event procedures registered in the connection are copied

into the `PGresult`. (They do not get `PGEVT_RESULTCREATE` calls, but see [PQfireResultCreateEvents](#).) Note that [PQclear](#) should eventually be called on the object, just as with a `PGresult` returned by libpq itself.

`PQfireResultCreateEvents`

Fires a `PGEVT_RESULTCREATE` event (see [Section 35.14](#)) for each event procedure registered in the `PGresult` object. Returns non-zero for success, zero if any event procedure fails.

```
int PQfireResultCreateEvents(PGconn *conn, PGresult *res);
```

The `conn` argument is passed through to event procedures but not used directly. It can be `NULL` if the event procedures won't use it.

Event procedures that have already received a `PGEVT_RESULTCREATE` or `PGEVT_RESULTCOPY` event for this object are not fired again.

The main reason that this function is separate from [PQmakeEmptyPGresult](#) is that it is often appropriate to create a `PGresult` and fill it with data before invoking the event procedures.

`PQcopyResult`

Makes a copy of a `PGresult` object. The copy is not linked to the source result in any way and [PQclear](#) must be called when the copy is no longer needed. If the function fails, `NULL` is returned.

```
PGresult *PQcopyResult(const PGresult *src, int flags);
```

This is not intended to make an exact copy. The returned result is always put into `PGRES_TUPLES_OK` status, and does not copy any error message in the source. (It does copy the command status string, however.) The `flags` argument determines what else is copied. It is a bitwise OR of several flags. `PG_COPYRES_ATTRS` specifies copying the source result's attributes (column definitions). `PG_COPYRES_TUPLES` specifies copying the source result's tuples. (This implies copying the attributes, too.) `PG_COPYRES_NOTICEHOOKS` specifies copying the source result's notify hooks. `PG_COPYRES_EVENTS` specifies copying the source result's events. (But any instance data associated with the source is not copied.) The event procedures receive `PGEVT_RESULTCOPY` events.

`PQsetResultAttrs`

Sets the attributes of a `PGresult` object.

```
int PQsetResultAttrs(PGresult *res, int numAttributes, PGresAttDesc *attDescs);
```

The provided `attDescs` are copied into the result. If the `attDescs` pointer is `NULL` or `numAttributes` is less than one, the request is ignored and the function succeeds. If `res` already contains attributes, the function will fail. If the function fails, the return value is zero. If the function succeeds, the return value is non-zero.

`PQsetvalue`

Sets a tuple field value of a `PGresult` object.

```
int PQsetvalue(PGresult *res, int tup_num, int field_num, char *value, int len);
```

The function will automatically grow the result's internal tuples array as needed. However, the `tup_num` argument must be less than or equal to [PQntuples](#), meaning this function can only grow the tuples array one tuple at a time. But any field of any existing tuple can be modified in any order. If a value at `field_num` already exists, it will be overwritten. If `len` is -1 or `value` is `NULL`, the field value will be set to an SQL null value. The `value` is copied into the result's private storage, thus is no longer needed after the function returns. If the function fails, the return value is zero. If the function succeeds, the return value is non-zero.

`PQresultAlloc`

Allocate subsidiary storage for a `PGresult` object.

```
void *PQresultAlloc(PGresult *res, size_t nBytes);
```

Any memory allocated with this function will be freed when *res* is cleared. If the function fails, the return value is `NULL`. The result is guaranteed to be adequately aligned for any type of data, just as for `malloc`.

`PQresultMemorySize`

Retrieves the number of bytes allocated for a `PGresult` object.

```
size_t PQresultMemorySize(const PGresult *res);
```

This value is the sum of all `malloc` requests associated with the `PGresult` object, that is, all the memory that will be freed by `PQclear`. This information can be useful for managing memory consumption.

`PQlibVersion`

Return the version of libpq that is being used.

```
int PQlibVersion(void);
```

The result of this function can be used to determine, at run time, whether specific functionality is available in the currently loaded version of libpq. The function can be used, for example, to determine which connection options are available in `PQconnectdb`.

The result is formed by multiplying the library's major version number by 10000 and adding the minor version number. For example, version 10.1 will be returned as 100001, and version 11.0 will be returned as 110000.

Prior to major version 10, Postgres Pro used three-part version numbers in which the first two parts together represented the major version. For those versions, `PQlibVersion` uses two digits for each part; for example version 9.1.5 will be returned as 90105, and version 9.2.0 will be returned as 90200.

Therefore, for purposes of determining feature compatibility, applications should divide the result of `PQlibVersion` by 100 not 10000 to determine a logical major version number. In all release series, only the last two digits differ between minor releases (bug-fix releases).

Note

This function appeared in PostgreSQL version 9.1, so it cannot be used to detect required functionality in earlier versions, since calling it will create a link dependency on version 9.1 or later.

`PQgetcurrentTimeUsec`

Retrieves the current time, expressed as the number of microseconds since the Unix epoch (that is, `time_t` times 1 million).

```
pg_usec_time_t PQgetcurrentTimeUsec(void);
```

This is primarily useful for calculating timeout values to use with `PQsocketPoll`.

35.13. Notice Processing

Notice and warning messages generated by the server are not returned by the query execution functions, since they do not imply failure of the query. Instead they are passed to a notice handling function, and execution continues normally after the handler returns. The default notice handling function prints the message on `stderr`, but the application can override this behavior by supplying its own handling function.

For historical reasons, there are two levels of notice handling, called the notice receiver and notice processor. The default behavior is for the notice receiver to format the notice and pass a string to the

notice processor for printing. However, an application that chooses to provide its own notice receiver will typically ignore the notice processor layer and just do all the work in the notice receiver.

The function `PQsetNoticeReceiver` sets or examines the current notice receiver for a connection object. Similarly, `PQsetNoticeProcessor` sets or examines the current notice processor.

```
typedef void (*PQnoticeReceiver) (void *arg, const PGresult *res);

PQnoticeReceiver
PQsetNoticeReceiver(PGconn *conn,
                   PQnoticeReceiver proc,
                   void *arg);

typedef void (*PQnoticeProcessor) (void *arg, const char *message);

PQnoticeProcessor
PQsetNoticeProcessor(PGconn *conn,
                   PQnoticeProcessor proc,
                   void *arg);
```

Each of these functions returns the previous notice receiver or processor function pointer, and sets the new value. If you supply a null function pointer, no action is taken, but the current pointer is returned.

When a notice or warning message is received from the server, or generated internally by libpq, the notice receiver function is called. It is passed the message in the form of a `PGRES_NONFATAL_ERROR` `PGresult`. (This allows the receiver to extract individual fields using `PQresultErrorField`, or obtain a complete preformatted message using `PQresultErrorMessage` or `PQresultVerboseErrorMessage`.) The same void pointer passed to `PQsetNoticeReceiver` is also passed. (This pointer can be used to access application-specific state if needed.)

The default notice receiver simply extracts the message (using `PQresultErrorMessage`) and passes it to the notice processor.

The notice processor is responsible for handling a notice or warning message given in text form. It is passed the string text of the message (including a trailing newline), plus a void pointer that is the same one passed to `PQsetNoticeProcessor`. (This pointer can be used to access application-specific state if needed.)

The default notice processor is simply:

```
static void
defaultNoticeProcessor(void *arg, const char *message)
{
    fprintf(stderr, "%s", message);
}
```

Once you have set a notice receiver or processor, you should expect that that function could be called as long as either the `PGconn` object or `PGresult` objects made from it exist. At creation of a `PGresult`, the `PGconn`'s current notice handling pointers are copied into the `PGresult` for possible use by functions like `PQgetvalue`.

35.14. Event System

libpq's event system is designed to notify registered event handlers about interesting libpq events, such as the creation or destruction of `PGconn` and `PGresult` objects. A principal use case is that this allows applications to associate their own data with a `PGconn` or `PGresult` and ensure that that data is freed at an appropriate time.

Each registered event handler is associated with two pieces of data, known to libpq only as opaque `void *` pointers. There is a *pass-through* pointer that is provided by the application when the event handler

is registered with a `PGconn`. The pass-through pointer never changes for the life of the `PGconn` and all `PGresult`s generated from it; so if used, it must point to long-lived data. In addition there is an *instance data* pointer, which starts out `NULL` in every `PGconn` and `PGresult`. This pointer can be manipulated using the `PQinstanceData`, `PQsetInstanceData`, `PQresultInstanceData` and `PQresultSetInstanceData` functions. Note that unlike the pass-through pointer, instance data of a `PGconn` is not automatically inherited by `PGresult`s created from it. libpq does not know what pass-through and instance data pointers point to (if anything) and will never attempt to free them — that is the responsibility of the event handler.

35.14.1. Event Types

The enum `PGEventId` names the types of events handled by the event system. All its values have names beginning with `PGEVT`. For each event type, there is a corresponding event info structure that carries the parameters passed to the event handlers. The event types are:

`PGEVT_REGISTER`

The register event occurs when `PQregisterEventProc` is called. It is the ideal time to initialize any `instanceData` an event procedure may need. Only one register event will be fired per event handler per connection. If the event procedure fails (returns zero), the registration is canceled.

```
typedef struct
{
    PGconn *conn;
} PGEventRegister;
```

When a `PGEVT_REGISTER` event is received, the `evtInfo` pointer should be cast to a `PGEventRegister *`. This structure contains a `PGconn` that should be in the `CONNECTION_OK` status; guaranteed if one calls `PQregisterEventProc` right after obtaining a good `PGconn`. When returning a failure code, all cleanup must be performed as no `PGEVT_CONNDESTROY` event will be sent.

`PGEVT_CONNRESET`

The connection reset event is fired on completion of `PQreset` or `PQresetPoll`. In both cases, the event is only fired if the reset was successful. The return value of the event procedure is ignored in Postgres Pro v15 and later. With earlier versions, however, it's important to return success (nonzero) or the connection will be aborted.

```
typedef struct
{
    PGconn *conn;
} PGEventConnReset;
```

When a `PGEVT_CONNRESET` event is received, the `evtInfo` pointer should be cast to a `PGEventConnReset *`. Although the contained `PGconn` was just reset, all event data remains unchanged. This event should be used to reset/reload/requery any associated `instanceData`. Note that even if the event procedure fails to process `PGEVT_CONNRESET`, it will still receive a `PGEVT_CONNDESTROY` event when the connection is closed.

`PGEVT_CONNDESTROY`

The connection destroy event is fired in response to `PQfinish`. It is the event procedure's responsibility to properly clean up its event data as libpq has no ability to manage this memory. Failure to clean up will lead to memory leaks.

```
typedef struct
{
    PGconn *conn;
} PGEventConnDestroy;
```

When a `PGEVT_CONNDESTROY` event is received, the `evtInfo` pointer should be cast to a `PGEventConnDestroy *`. This event is fired prior to `PQfinish` performing any other cleanup. The return value of the event procedure is ignored since there is no way of indicating a failure from `PQfinish`. Also, an event procedure failure should not abort the process of cleaning up unwanted memory.

PGEVT_RESULTCREATE

The result creation event is fired in response to any query execution function that generates a result, including [PQgetResult](#). This event will only be fired after the result has been created successfully.

```
typedef struct
{
    PGconn *conn;
    PGresult *result;
} PGEventResultCreate;
```

When a `PGEVT_RESULTCREATE` event is received, the *evtInfo* pointer should be cast to a `PGEventResultCreate *`. The *conn* is the connection used to generate the result. This is the ideal place to initialize any *instanceData* that needs to be associated with the result. If an event procedure fails (returns zero), that event procedure will be ignored for the remaining lifetime of the result; that is, it will not receive `PGEVT_RESULTCOPY` or `PGEVT_RESULTDESTROY` events for this result or results copied from it.

PGEVT_RESULTCOPY

The result copy event is fired in response to [PQcopyResult](#). This event will only be fired after the copy is complete. Only event procedures that have successfully handled the `PGEVT_RESULTCREATE` or `PGEVT_RESULTCOPY` event for the source result will receive `PGEVT_RESULTCOPY` events.

```
typedef struct
{
    const PGresult *src;
    PGresult *dest;
} PGEventResultCopy;
```

When a `PGEVT_RESULTCOPY` event is received, the *evtInfo* pointer should be cast to a `PGEventResultCopy *`. The *src* result is what was copied while the *dest* result is the copy destination. This event can be used to provide a deep copy of *instanceData*, since `PQcopyResult` cannot do that. If an event procedure fails (returns zero), that event procedure will be ignored for the remaining lifetime of the new result; that is, it will not receive `PGEVT_RESULTCOPY` or `PGEVT_RESULTDESTROY` events for that result or results copied from it.

PGEVT_RESULTDESTROY

The result destroy event is fired in response to a [PQclear](#). It is the event procedure's responsibility to properly clean up its event data as libpq has no ability to manage this memory. Failure to clean up will lead to memory leaks.

```
typedef struct
{
    PGresult *result;
} PGEventResultDestroy;
```

When a `PGEVT_RESULTDESTROY` event is received, the *evtInfo* pointer should be cast to a `PGEventResultDestroy *`. This event is fired prior to [PQclear](#) performing any other cleanup. The return value of the event procedure is ignored since there is no way of indicating a failure from [PQclear](#). Also, an event procedure failure should not abort the process of cleaning up unwanted memory.

35.14.2. Event Callback Procedure

PGEventProc

`PGEventProc` is a typedef for a pointer to an event procedure, that is, the user callback function that receives events from libpq. The signature of an event procedure must be

```
int eventproc(PGEventId evtId, void *evtInfo, void *passThrough)
```

The *evtId* parameter indicates which `PGEVT` event occurred. The *evtInfo* pointer must be cast to the appropriate structure type to obtain further information about the event. The *passThrough* pa-

parameter is the pointer provided to `PQregisterEventProc` when the event procedure was registered. The function should return a non-zero value if it succeeds and zero if it fails.

A particular event procedure can be registered only once in any `PGconn`. This is because the address of the procedure is used as a lookup key to identify the associated instance data.

Caution

On Windows, functions can have two different addresses: one visible from outside a DLL and another visible from inside the DLL. One should be careful that only one of these addresses is used with libpq's event-procedure functions, else confusion will result. The simplest rule for writing code that will work is to ensure that event procedures are declared `static`. If the procedure's address must be available outside its own source file, expose a separate function to return the address.

35.14.3. Event Support Functions

`PQregisterEventProc`

Registers an event callback procedure with libpq.

```
int PQregisterEventProc(PGconn *conn, PGEventProc proc,
                       const char *name, void *passThrough);
```

An event procedure must be registered once on each `PGconn` you want to receive events about. There is no limit, other than memory, on the number of event procedures that can be registered with a connection. The function returns a non-zero value if it succeeds and zero if it fails.

The *proc* argument will be called when a libpq event is fired. Its memory address is also used to lookup *instanceData*. The *name* argument is used to refer to the event procedure in error messages. This value cannot be `NULL` or a zero-length string. The name string is copied into the `PGconn`, so what is passed need not be long-lived. The *passThrough* pointer is passed to the *proc* whenever an event occurs. This argument can be `NULL`.

`PQsetInstanceData`

Sets the connection *conn*'s *instanceData* for procedure *proc* to *data*. This returns non-zero for success and zero for failure. (Failure is only possible if *proc* has not been properly registered in *conn*.)

```
int PQsetInstanceData(PGconn *conn, PGEventProc proc, void *data);
```

`PQinstanceData`

Returns the connection *conn*'s *instanceData* associated with procedure *proc*, or `NULL` if there is none.

```
void *PQinstanceData(const PGconn *conn, PGEventProc proc);
```

`PQresultSetInstanceData`

Sets the result's *instanceData* for *proc* to *data*. This returns non-zero for success and zero for failure. (Failure is only possible if *proc* has not been properly registered in the result.)

```
int PQresultSetInstanceData(PGresult *res, PGEventProc proc, void *data);
```

Beware that any storage represented by *data* will not be accounted for by `PQresultMemorySize`, unless it is allocated using `PQresultAlloc`. (Doing so is recommendable because it eliminates the need to free such storage explicitly when the result is destroyed.)

`PQresultInstanceData`

Returns the result's *instanceData* associated with *proc*, or `NULL` if there is none.

```
void *PQresultInstanceData(const PGresult *res, PGEventProc proc);
```


35.14.4. Event Example

Here is a skeleton example of managing private data associated with libpq connections and results.

```
/* required header for libpq events (note: includes libpq-fe.h) */
#include <libpq-events.h>

/* The instanceData */
typedef struct
{
    int n;
    char *str;
} mydata;

/* PGEventProc */
static int myEventProc(PGEventId evtId, void *evtInfo, void *passThrough);

int
main(void)
{
    mydata *data;
    PGresult *res;
    PGconn *conn =
        PQconnectdb("dbname=postgres options=-csearch_path=");

    if (PQstatus(conn) != CONNECTION_OK)
    {
        /* PQerrorMessage's result includes a trailing newline */
        fprintf(stderr, "%s", PQerrorMessage(conn));
        PQfinish(conn);
        return 1;
    }

    /* called once on any connection that should receive events.
     * Sends a PGEVT_REGISTER to myEventProc.
     */
    if (!PQregisterEventProc(conn, myEventProc, "mydata_proc", NULL))
    {
        fprintf(stderr, "Cannot register PGEventProc\n");
        PQfinish(conn);
        return 1;
    }

    /* conn instanceData is available */
    data = PQinstanceData(conn, myEventProc);

    /* Sends a PGEVT_RESULTCREATE to myEventProc */
    res = PQexec(conn, "SELECT 1 + 1");

    /* result instanceData is available */
    data = PQresultInstanceData(res, myEventProc);

    /* If PG_COPYRES_EVENTS is used, sends a PGEVT_RESULTCOPY to myEventProc */
    res_copy = PQcopyResult(res, PG_COPYRES_TUPLES | PG_COPYRES_EVENTS);

    /* result instanceData is available if PG_COPYRES_EVENTS was
     * used during the PQcopyResult call.
     */
}
```



```
    */
    data = PQresultInstanceData(res_copy, myEventProc);

    /* Both clears send a PGEVT_RESULTDESTROY to myEventProc */
    PQclear(res);
    PQclear(res_copy);

    /* Sends a PGEVT_CONNDESTROY to myEventProc */
    PQfinish(conn);

    return 0;
}

static int
myEventProc(PGEventId evtId, void *evtInfo, void *passThrough)
{
    switch (evtId)
    {
        case PGEVT_REGISTER:
        {
            PGEventRegister *e = (PGEventRegister *)evtInfo;
            mydata *data = get_mydata(e->conn);

            /* associate app specific data with connection */
            PQsetInstanceData(e->conn, myEventProc, data);
            break;
        }

        case PGEVT_CONNRESET:
        {
            PGEventConnReset *e = (PGEventConnReset *)evtInfo;
            mydata *data = PQinstanceData(e->conn, myEventProc);

            if (data)
                memset(data, 0, sizeof(mydata));
            break;
        }

        case PGEVT_CONNDESTROY:
        {
            PGEventConnDestroy *e = (PGEventConnDestroy *)evtInfo;
            mydata *data = PQinstanceData(e->conn, myEventProc);

            /* free instance data because the conn is being destroyed */
            if (data)
                free_mydata(data);
            break;
        }

        case PGEVT_RESULTCREATE:
        {
            PGEventResultCreate *e = (PGEventResultCreate *)evtInfo;
            mydata *conn_data = PQinstanceData(e->conn, myEventProc);
            mydata *res_data = dup_mydata(conn_data);

            /* associate app specific data with result (copy it from conn) */
            PQresultSetInstanceData(e->result, myEventProc, res_data);
            break;
        }
    }
}
```

```
    }

    case PGEVT_RESULTCOPY:
    {
        PGEventResultCopy *e = (PGEventResultCopy *)evtInfo;
        mydata *src_data = PQresultInstanceData(e->src, myEventProc);
        mydata *dest_data = dup_mydata(src_data);

        /* associate app specific data with result (copy it from a result) */
        PQresultSetInstanceData(e->dest, myEventProc, dest_data);
        break;
    }

    case PGEVT_RESULTDESTROY:
    {
        PGEventResultDestroy *e = (PGEventResultDestroy *)evtInfo;
        mydata *data = PQresultInstanceData(e->result, myEventProc);

        /* free instance data because the result is being destroyed */
        if (data)
            free_mydata(data);
        break;
    }

    /* unknown event ID, just return true. */
    default:
        break;
}

return true; /* event processing succeeded */
}
```

35.15. Environment Variables

The following environment variables can be used to select default connection parameter values, which will be used by [PQconnectdb](#), [PQsetdbLogin](#) and [PQsetdb](#) if no value is directly specified by the calling code. These are useful to avoid hard-coding database connection information into simple client applications, for example.

- `PGHOST` behaves the same as the [host](#) connection parameter.
- `PGSSLNEGOTIATION` behaves the same as the [sslnegotiation](#) connection parameter.
- `PGHOSTADDR` behaves the same as the [hostaddr](#) connection parameter. This can be set instead of or in addition to `PGHOST` to avoid DNS lookup overhead.
- `PGPORT` behaves the same as the [port](#) connection parameter.
- `PGDATABASE` behaves the same as the [dbname](#) connection parameter.
- `PGUSER` behaves the same as the [user](#) connection parameter.
- `PGPASSWORD` behaves the same as the [password](#) connection parameter. Use of this environment variable is not recommended for security reasons, as some operating systems allow non-root users to see process environment variables via `ps`; instead consider using a password file (see [Section 35.16](#)).
- `PGPASSFILE` behaves the same as the [passfile](#) connection parameter.
- `PGREQUIREAUTH` behaves the same as the [require_auth](#) connection parameter.
- `PGCHANNELBINDING` behaves the same as the [channel_binding](#) connection parameter.
- `PGSERVICE` behaves the same as the [service](#) connection parameter.

- `PGSERVICEFILE` specifies the name of the per-user connection service file (see [Section 35.17](#)). Defaults to `~/.pg_service.conf`, or `%APPDATA%\postgresql\.pg_service.conf` on Microsoft Windows.
- `PGOPTIONS` behaves the same as the [options](#) connection parameter.
- `PGAPPNAME` behaves the same as the [application_name](#) connection parameter.
- `PGSSLMODE` behaves the same as the [sslmode](#) connection parameter.
- `PGREQUIRESSL` behaves the same as the [requiressl](#) connection parameter. This environment variable is deprecated in favor of the `PGSSLMODE` variable; setting both variables suppresses the effect of this one.
- `PGSSLCOMPRESSION` behaves the same as the [sslcompression](#) connection parameter.
- `PGSSLCERT` behaves the same as the [sslcert](#) connection parameter.
- `PGSSLKEY` behaves the same as the [sslkey](#) connection parameter.
- `PGSSLCERTMODE` behaves the same as the [sslcertmode](#) connection parameter.
- `PGSSLROOTCERT` behaves the same as the [sslrootcert](#) connection parameter.
- `PGSSLCRL` behaves the same as the [sslcrl](#) connection parameter.
- `PGSSLCRLDIR` behaves the same as the [sslcrlidir](#) connection parameter.
- `PGSSLSNI` behaves the same as the [sslsni](#) connection parameter.
- `PGREQUIREPEER` behaves the same as the [requirepeer](#) connection parameter.
- `PGSSLMINPROTOCOLVERSION` behaves the same as the [ssl_min_protocol_version](#) connection parameter.
- `PGSSLMAXPROTOCOLVERSION` behaves the same as the [ssl_max_protocol_version](#) connection parameter.
- `PGGSSENCMODE` behaves the same as the [gssencmode](#) connection parameter.
- `PGKRBSRVNAME` behaves the same as the [krbsrvname](#) connection parameter.
- `PGGSSLIB` behaves the same as the [gsslib](#) connection parameter.
- `PGGSSDELEGATION` behaves the same as the [gssdelegation](#) connection parameter.
- `PGCONNECT_TIMEOUT` behaves the same as the [connect_timeout](#) connection parameter.
- `PGCLIENTENCODING` behaves the same as the [client_encoding](#) connection parameter.
- `PGTARGETSESSIONATTRS` behaves the same as the [target_session_attrs](#) connection parameter.
- `PGLOADBALANCEHOSTS` behaves the same as the [load_balance_hosts](#) connection parameter.

The following environment variables can be used to specify default behavior for each Postgres Pro session. (See also the [ALTER ROLE](#) and [ALTER DATABASE](#) commands for ways to set default behavior on a per-user or per-database basis.)

- `PGDATESTYLE` sets the default style of date/time representation. (Equivalent to `SET datestyle TO ...`.)
- `PGTZ` sets the default time zone. (Equivalent to `SET timezone TO ...`.)
- `PGGEQO` sets the default mode for the genetic query optimizer. (Equivalent to `SET geqo TO ...`.)

Refer to the SQL command [SET](#) for information on correct values for these environment variables.

The following environment variables determine internal behavior of libpq; they override compiled-in defaults.

- `PGSYSCONFDIR` sets the directory containing the `pg_service.conf` file and in a future version possibly other system-wide configuration files.
- `PGREUSEPASS` behaves the same as the [reusepass](#) connection parameter. Use of this environment variable can disable automatic reconnection.
- `PGLOCALEDIR` sets the directory containing the `locale` files for message localization.

35.16. The Password File

The file `.pgpass` in a user's home directory can contain passwords to be used if the connection requires a password (and no password has been specified otherwise). On Microsoft Windows the file is named `%APPDATA%\postgresql\pgpass.conf` (where `%APPDATA%` refers to the Application Data subdirectory in the user's profile). Alternatively, the password file to use can be specified using the connection parameter `passfile` or the environment variable `PGPASSFILE`.

This file should contain lines of the following format:

```
hostname:port:database:username:password
```

(You can add a reminder comment to the file by copying the line above and preceding it with `#`.) Each of the first four fields can be a literal value, or `*`, which matches anything. The password field from the first line that matches the current connection parameters will be used. (Therefore, put more-specific entries first when you are using wildcards.) If an entry needs to contain `:` or `\`, escape this character with `\`. The host name field is matched to the `host` connection parameter if that is specified, otherwise to the `hostaddr` parameter if that is specified; if neither are given then the host name `localhost` is searched for. The host name `localhost` is also searched for when the connection is a Unix-domain socket connection and the `host` parameter matches libpq's default socket directory path. In a standby server, a database field of `replication` matches streaming replication connections made to the primary server. The database field is of limited usefulness otherwise, because users have the same password for all databases in the same cluster.

On Unix systems, the permissions on a password file must disallow any access to world or group; achieve this by a command such as `chmod 0600 ~/.pgpass`. If the permissions are less strict than this, the file will be ignored. On Microsoft Windows, it is assumed that the file is stored in a directory that is secure, so no special permissions check is made.

35.17. The Connection Service File

The connection service file allows libpq connection parameters to be associated with a single service name. That service name can then be specified in a libpq connection string, and the associated settings will be used. This allows connection parameters to be modified without requiring a recompile of the libpq-using application. The service name can also be specified using the `PGSERVICE` environment variable.

Service names can be defined in either a per-user service file or a system-wide file. If the same service name exists in both the user and the system file, the user file takes precedence. By default, the per-user service file is named `~/.pg_service.conf`. On Microsoft Windows, it is named `%APPDATA%\postgresql\pg_service.conf` (where `%APPDATA%` refers to the Application Data subdirectory in the user's profile). A different file name can be specified by setting the environment variable `PGSERVICEFILE`. The system-wide file is named `pg_service.conf`. By default it is sought in the `etc` directory of the Postgres Pro installation (use `pg_config --sysconfdir` to identify this directory precisely). Another directory, but not a different file name, can be specified by setting the environment variable `PGSYSCONFDIR`.

Either service file uses an “INI file” format where the section name is the service name and the parameters are connection parameters; see [Section 35.1.2](#) for a list. For example:

```
# comment
[mydb]
host=somehost
port=5433
user=admin
```

An example file is provided in the Postgres Pro installation at `share/pg_service.conf.sample`.

Connection parameters obtained from a service file are combined with parameters obtained from other sources. A service file setting overrides the corresponding environment variable, and in turn can be overridden by a value given directly in the connection string. For example, using the above service file, a connection string `service=mydb port=5434` will use host `somehost`, port `5434`, user `admin`, and other parameters as set by environment variables or built-in defaults.

35.18. LDAP Lookup of Connection Parameters

If libpq has been compiled with LDAP support (option `--with-ldap` for configure) it is possible to retrieve connection options like `host` or `dbname` via LDAP from a central server. The advantage is that if the connection parameters for a database change, the connection information doesn't have to be updated on all client machines.

LDAP connection parameter lookup uses the connection service file `pg_service.conf` (see [Section 35.17](#)). A line in a `pg_service.conf` stanza that starts with `ldap://` will be recognized as an LDAP URL and an LDAP query will be performed. The result must be a list of `keyword = value` pairs which will be used to set connection options. The URL must conform to [RFC 1959](#) and be of the form

```
ldap://[hostname[:port]]/search_base?attribute?search_scope?filter
```

where `hostname` defaults to `localhost` and `port` defaults to `389`.

Processing of `pg_service.conf` is terminated after a successful LDAP lookup, but is continued if the LDAP server cannot be contacted. This is to provide a fallback with further LDAP URL lines that point to different LDAP servers, classical `keyword = value` pairs, or default connection options. If you would rather get an error message in this case, add a syntactically incorrect line after the LDAP URL.

A sample LDAP entry that has been created with the LDIF file

```
version:1
dn:cn=mydatabase,dc=mycompany,dc=com
changetype:add
objectclass:top
objectclass:device
cn:mydatabase
description:host=dbserver.mycompany.com
description:port=5439
description:dbname=mydb
description:user=mydb_user
description:sslmode=require
```

might be queried with the following LDAP URL:

```
ldap://ldap.mycompany.com/dc=mycompany,dc=com?description?one?(cn=mydatabase)
```

You can also mix regular service file entries with LDAP lookups. A complete example for a stanza in `pg_service.conf` would be:

```
# only host and port are stored in LDAP, specify dbname and user explicitly
[customerdb]
dbname=customer
user=appuser
ldap://ldap.acme.com/cn=dbserver,cn=hosts?pgconnectinfo?base?(objectclass=*)
```

35.19. SSL Support

Postgres Pro has native support for using SSL connections to encrypt client/server communications using TLS protocols for increased security. See [Section 18.9](#) for details about the server-side SSL functionality.

libpq reads the system-wide OpenSSL configuration file. By default, this file is named `openssl.cnf` and is located in the directory reported by `openssl version -d`. This default can be overridden by setting environment variable `OPENSSL_CONF` to the name of the desired configuration file.

35.19.1. Client Verification of Server Certificates

By default, Postgres Pro will not perform any verification of the server certificate. This means that it is possible to spoof the server identity (for example by modifying a DNS record or by taking over the server IP address) without the client knowing. In order to prevent spoofing, the client must be able to verify the server's identity via a chain of trust. A chain of trust is established by placing a root (self-

signed) certificate authority (CA) certificate on one computer and a leaf certificate *signed* by the root certificate on another computer. It is also possible to use an “intermediate” certificate which is signed by the root certificate and signs leaf certificates.

To allow the client to verify the identity of the server, place a root certificate on the client and a leaf certificate signed by the root certificate on the server. To allow the server to verify the identity of the client, place a root certificate on the server and a leaf certificate signed by the root certificate on the client. One or more intermediate certificates (usually stored with the leaf certificate) can also be used to link the leaf certificate to the root certificate.

Once a chain of trust has been established, there are two ways for the client to validate the leaf certificate sent by the server. If the parameter `sslmode` is set to `verify-ca`, libpq will verify that the server is trustworthy by checking the certificate chain up to the root certificate stored on the client. If `sslmode` is set to `verify-full`, libpq will *also* verify that the server host name matches the name stored in the server certificate. The SSL connection will fail if the server certificate cannot be verified. `verify-full` is recommended in most security-sensitive environments.

In `verify-full` mode, the host name is matched against the certificate's Subject Alternative Name attribute(s) (SAN), or against the Common Name attribute if no SAN of type `dNSName` is present. If the certificate's name attribute starts with an asterisk (*), the asterisk will be treated as a wildcard, which will match all characters *except* a dot (.). This means the certificate will not match subdomains. If the connection is made using an IP address instead of a host name, the IP address will be matched (without doing any DNS lookups) against SANs of type `iPAddress` or `dNSName`. If no `iPAddress` SAN is present and no matching `dNSName` SAN is present, the host IP address is matched against the Common Name attribute.

Note

For backward compatibility with earlier versions of Postgres Pro, the host IP address is verified in a manner different from [RFC 6125](#). The host IP address is always matched against `dNSName` SANs as well as `iPAddress` SANs, and can be matched against the Common Name attribute if no relevant SANs exist.

To allow server certificate verification, one or more root certificates must be placed in the file `~/.postgresql/root.crt` in the user's home directory. (On Microsoft Windows the file is named `%APPDATA%\postgresql\root.crt`.) Intermediate certificates should also be added to the file if they are needed to link the certificate chain sent by the server to the root certificates stored on the client.

Certificate Revocation List (CRL) entries are also checked if the file `~/.postgresql/root.crl` exists (`%APPDATA%\postgresql\root.crl` on Microsoft Windows).

The location of the root certificate file and the CRL can be changed by setting the connection parameters `sslrootcert` and `sslcrldir` or the environment variables `PGSSLROOTCERT` and `PGSSLCRL`. `sslcrldir` or the environment variable `PGSSLCRLDIR` can also be used to specify a directory containing CRL files.

Note

For backwards compatibility with earlier versions of Postgres Pro, if a root CA file exists, the behavior of `sslmode=require` will be the same as that of `verify-ca`, meaning the server certificate is validated against the CA. Relying on this behavior is discouraged, and applications that need certificate validation should always use `verify-ca` or `verify-full`.

35.19.2. Client Certificates

If the server attempts to verify the identity of the client by requesting the client's leaf certificate, libpq will send the certificate(s) stored in file `~/.postgresql/postgresql.crt` in the user's home directory. The certificates must chain to the root certificate trusted by the server. A matching private key file

`~/.postgresql/postgresql.key` must also be present. On Microsoft Windows these files are named `%APPDATA%\postgresql\postgresql.crt` and `%APPDATA%\postgresql\postgresql.key`. The location of the certificate and key files can be overridden by the connection parameters `sslcert` and `sslkey`, or by the environment variables `PGSSLCERT` and `PGSSLKEY`.

On Unix systems, the permissions on the private key file must disallow any access to world or group; achieve this by a command such as `chmod 0600 ~/.postgresql/postgresql.key`. Alternatively, the file can be owned by root and have group read access (that is, 0640 permissions). That setup is intended for installations where certificate and key files are managed by the operating system. The user of libpq should then be made a member of the group that has access to those certificate and key files. (On Microsoft Windows, there is no file permissions check, since the `%APPDATA%\postgresql` directory is presumed secure.)

The first certificate in `postgresql.crt` must be the client's certificate because it must match the client's private key. “Intermediate” certificates can be optionally appended to the file — doing so avoids requiring storage of intermediate certificates on the server ([ssl_ca_file](#)).

The certificate and key may be in PEM or ASN.1 DER format.

The key may be stored in cleartext or encrypted with a passphrase using any algorithm supported by OpenSSL, like AES-128. If the key is stored encrypted, then the passphrase may be provided in the [sslpassword](#) connection option. If an encrypted key is supplied and the `sslpassword` option is absent or blank, a password will be prompted for interactively by OpenSSL with a `Enter PEM passphrase: prompt` if a TTY is available. Applications can override the client certificate prompt and the handling of the `sslpassword` parameter by supplying their own key password callback; see [PQsetSSLKey-PassHook_OpenSSL](#).

For instructions on creating certificates, see [Section 18.9.5](#).

35.19.3. Protection Provided in Different Modes

The different values for the `sslmode` parameter provide different levels of protection. SSL can provide protection against three types of attacks:

Eavesdropping

If a third party can examine the network traffic between the client and the server, it can read both connection information (including the user name and password) and the data that is passed. SSL uses encryption to prevent this.

Man-in-the-middle (MITM)

If a third party can modify the data while passing between the client and server, it can pretend to be the server and therefore see and modify data *even if it is encrypted*. The third party can then forward the connection information and data to the original server, making it impossible to detect this attack. Common vectors to do this include DNS poisoning and address hijacking, whereby the client is directed to a different server than intended. There are also several other attack methods that can accomplish this. SSL uses certificate verification to prevent this, by authenticating the server to the client.

Impersonation

If a third party can pretend to be an authorized client, it can simply access data it should not have access to. Typically this can happen through insecure password management. SSL uses client certificates to prevent this, by making sure that only holders of valid certificates can access the server.

For a connection to be known SSL-secured, SSL usage must be configured on *both the client and the server* before the connection is made. If it is only configured on the server, the client may end up sending sensitive information (e.g., passwords) before it knows that the server requires high security. In libpq, secure connections can be ensured by setting the `sslmode` parameter to `verify-full` or `verify-ca`, and providing the system with a root certificate to verify against. This is analogous to using an `https` URL for encrypted web browsing.

Once the server has been authenticated, the client can pass sensitive data. This means that up until this point, the client does not need to know if certificates will be used for authentication, making it safe to specify that only in the server configuration.

All SSL options carry overhead in the form of encryption and key-exchange, so there is a trade-off that has to be made between performance and security. [Table 35.1](#) illustrates the risks the different `sslmode` values protect against, and what statement they make about security and overhead.

Table 35.1. SSL Mode Descriptions

<code>sslmode</code>	Eavesdropping protection	MITM protection	Statement
<code>disable</code>	No	No	I don't care about security, and I don't want to pay the overhead of encryption.
<code>allow</code>	Maybe	No	I don't care about security, but I will pay the overhead of encryption if the server insists on it.
<code>prefer</code>	Maybe	No	I don't care about encryption, but I wish to pay the overhead of encryption if the server supports it.
<code>require</code>	Yes	No	I want my data to be encrypted, and I accept the overhead. I trust that the network will make sure I always connect to the server I want.
<code>verify-ca</code>	Yes	Depends on CA policy	I want my data encrypted, and I accept the overhead. I want to be sure that I connect to a server that I trust.
<code>verify-full</code>	Yes	Yes	I want my data encrypted, and I accept the overhead. I want to be sure that I connect to a server I trust, and that it's the one I specify.

The difference between `verify-ca` and `verify-full` depends on the policy of the root CA. If a public CA is used, `verify-ca` allows connections to a server that *somebody else* may have registered with the CA. In this case, `verify-full` should always be used. If a local CA is used, or even a self-signed certificate, using `verify-ca` often provides enough protection.

The default value for `sslmode` is `prefer`. As is shown in the table, this makes no sense from a security point of view, and it only promises performance overhead if possible. It is only provided as the default for backward compatibility, and is not recommended in secure deployments.

35.19.4. SSL Client File Usage

[Table 35.2](#) summarizes the files that are relevant to the SSL setup on the client.

Table 35.2. Libpq/Client SSL File Usage

File	Contents	Effect
<code>~/.postgresql/postgresql.crt</code>	client certificate	sent to server
<code>~/.postgresql/postgresql.key</code>	client private key	proves client certificate sent by owner; does not indicate certificate owner is trustworthy
<code>~/.postgresql/root.crt</code>	trusted certificate authorities	checks that server certificate is signed by a trusted certificate authority
<code>~/.postgresql/root.crl</code>	certificates revoked by certificate authorities	server certificate must not be on this list

35.19.5. SSL Library Initialization

If your application initializes `libssl` and/or `libcrypto` libraries and `libpq` is built with SSL support, you should call `PQinitOpenSSL` to tell `libpq` that the `libssl` and/or `libcrypto` libraries have been initialized by your application, so that `libpq` will not also initialize those libraries. However, this is unnecessary when using OpenSSL version 1.1.0 or later, as duplicate initializations are no longer problematic.

`PQinitOpenSSL`

Allows applications to select which security libraries to initialize.

```
void PQinitOpenSSL(int do_ssl, int do_crypto);
```

When `do_ssl` is non-zero, `libpq` will initialize the OpenSSL library before first opening a database connection. When `do_crypto` is non-zero, the `libcrypto` library will be initialized. By default (if `PQinitOpenSSL` is not called), both libraries are initialized. When SSL support is not compiled in, this function is present but does nothing.

If your application uses and initializes either OpenSSL or its underlying `libcrypto` library, you *must* call this function with zeroes for the appropriate parameter(s) before first opening a database connection. Also be sure that you have done that initialization before opening a database connection.

`PQinitSSL`

Allows applications to select which security libraries to initialize.

```
void PQinitSSL(int do_ssl);
```

This function is equivalent to `PQinitOpenSSL(do_ssl, do_ssl)`. It is sufficient for applications that initialize both or neither of OpenSSL and `libcrypto`.

`PQinitSSL` has been present since PostgreSQL 8.0, while `PQinitOpenSSL` was added in PostgreSQL 8.4, so `PQinitSSL` might be preferable for applications that need to work with older versions of `libpq`.

35.20. Behavior in Threaded Programs

As of version 17, `libpq` is always reentrant and thread-safe. However, one restriction is that no two threads attempt to manipulate the same `PGconn` object at the same time. In particular, you cannot issue concurrent commands from different threads through the same connection object. (If you need to run concurrent commands, use multiple connections.)

`PGresult` objects are normally read-only after creation, and so can be passed around freely between threads. However, if you use any of the `PGresult`-modifying functions described in [Section 35.12](#) or [Section 35.14](#), it's up to you to avoid concurrent operations on the same `PGresult`, too.

In earlier versions, `libpq` could be compiled with or without thread support, depending on compiler options. This function allows the querying of `libpq`'s thread-safe status:

`PQisthreadsafe`

Returns the thread safety status of the `libpq` library.

```
int PQisthreadsafe();
```

Returns 1 if the `libpq` is thread-safe and 0 if it is not. Always returns 1 on version 17 and above.

The deprecated functions `PQrequestCancel` and `PQoidStatus` are not thread-safe and should not be used in multithread programs. `PQrequestCancel` can be replaced by `PQcancelBlocking`. `PQoidStatus` can be replaced by `PQoidValue`.

If you are using Kerberos inside your application (in addition to inside `libpq`), you will need to do locking around Kerberos calls because Kerberos functions are not thread-safe.

35.21. Building libpq Programs

To build (i.e., compile and link) a program using `libpq` you need to do all of the following things:

- Include the `libpq-fe.h` header file:

```
#include <libpq-fe.h>
```

If you failed to do that then you will normally get error messages from your compiler similar to:

```
foo.c: In function `main':
foo.c:34: `PGconn' undeclared (first use in this function)
foo.c:35: `PGresult' undeclared (first use in this function)
foo.c:54: `CONNECTION_BAD' undeclared (first use in this function)
foo.c:68: `PGRES_COMMAND_OK' undeclared (first use in this function)
foo.c:95: `PGRES_TUPLES_OK' undeclared (first use in this function)
```

- Point your compiler to the directory where the Postgres Pro header files were installed, by supplying the `-I` option to your compiler. (In some cases the compiler will look into the directory in question by default, so you can omit this option.) For instance, your compile command line could look like:

```
cc -c -I/usr/local/pgsql/include testprog.c
```

If you are using makefiles then add the option to the `CPPFLAGS` variable:

```
CPPFLAGS += -I/usr/local/pgsql/include
```

If there is any chance that your program might be compiled by other users then you should not hardcode the directory location like that. Instead, you can run the utility `pg_config` to find out where the header files are on the local system:

```
$ pg_config --includedir
/usr/local/include
```

If you have `pkg-config` installed, you can run instead:

```
$ pkg-config --cflags libpq
-I/usr/local/include
```

Note that this will already include the `-I` in front of the path.

Failure to specify the correct option to the compiler will result in an error message such as:

```
testlibpq.c:8:22: libpq-fe.h: No such file or directory
```

- When linking the final program, specify the option `-lpq` so that the `libpq` library gets pulled in, as well as the option `-L` to point the compiler to the directory where the `libpq` library resides. (Again, the compiler will search some directories by default.) For maximum portability, put the `-L` option before the `-lpq` option. For example:

```
cc -o testprog testprog1.o testprog2.o -L/usr/local/pgsql/lib -lpq
```

You can find out the library directory using `pg_config` as well:

```
$ pg_config --libdir
/usr/local/pgsql/lib
```

Or again use `pkg-config`:

```
$ pkg-config --libs libpq
-L/usr/local/pgsql/lib -lpq
```

Note again that this prints the full options, not only the path.

Error messages that point to problems in this area could look like the following:

```
testlibpq.o: In function `main':
testlibpq.o(.text+0x60): undefined reference to `PQsetdbLogin'
testlibpq.o(.text+0x71): undefined reference to `PQstatus'
testlibpq.o(.text+0xa4): undefined reference to `PQerrorMessage'
```

This means you forgot `-lpq`.

`/usr/bin/ld: cannot find -lpq`

This means you forgot the `-L` option or did not specify the right directory.

35.22. Example Programs

These examples and others can be found in the directory `src/test/examples` in the PostgreSQL source code distribution.

Example 35.1. libpq Example Program 1

```
/*
 * src/test/examples/testlibpq.c
 *
 *
 * testlibpq.c
 *
 *      Test the C version of libpq, the PostgreSQL frontend library.
 */
#include <stdio.h>
#include <stdlib.h>
#include "libpq-fe.h"

static void
exit_nicely(PGconn *conn)
{
    PQfinish(conn);
    exit(1);
}

int
main(int argc, char **argv)
{
    const char *conninfo;
    PGconn      *conn;
    PGresult     *res;
    int          nFields;
    int          i,
                j;

    /*
     * If the user supplies a parameter on the command line, use it as the
     * conninfo string; otherwise default to setting dbname=postgres and using
     * environment variables or defaults for all other connection parameters.
     */
    if (argc > 1)
        conninfo = argv[1];
    else
        conninfo = "dbname = postgres";

    /* Make a connection to the database */
    conn = PQconnectdb(conninfo);

    /* Check to see that the backend connection was successfully made */
    if (PQstatus(conn) != CONNECTION_OK)
    {
```

```
    fprintf(stderr, "%s", PQerrorMessage(conn));
    exit_nicely(conn);
}

/* Set always-secure search path, so malicious users can't take control. */
res = PQexec(conn,
    "SELECT pg_catalog.set_config('search_path', '', false)");
if (PQresultStatus(res) != PGRES_TUPLES_OK)
{
    fprintf(stderr, "SET failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}

/*
 * Should PQclear PGresult whenever it is no longer needed to avoid memory
 * leaks
 */
PQclear(res);

/*
 * Our test case here involves using a cursor, for which we must be inside
 * a transaction block. We could do the whole thing with a single
 * PQexec() of "select * from pg_database", but that's too trivial to make
 * a good example.
 */

/* Start a transaction block */
res = PQexec(conn, "BEGIN");
if (PQresultStatus(res) != PGRES_COMMAND_OK)
{
    fprintf(stderr, "BEGIN command failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}
PQclear(res);

/*
 * Fetch rows from pg_database, the system catalog of databases
 */
res = PQexec(conn, "DECLARE myportal CURSOR FOR select * from pg_database");
if (PQresultStatus(res) != PGRES_COMMAND_OK)
{
    fprintf(stderr, "DECLARE CURSOR failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}
PQclear(res);

res = PQexec(conn, "FETCH ALL in myportal");
if (PQresultStatus(res) != PGRES_TUPLES_OK)
{
    fprintf(stderr, "FETCH ALL failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}

/* first, print out the attribute names */
```

```
nFields = PQnfields(res);
for (i = 0; i < nFields; i++)
    printf("%-15s", PQfname(res, i));
printf("\n\n");

/* next, print out the rows */
for (i = 0; i < PQntuples(res); i++)
{
    for (j = 0; j < nFields; j++)
        printf("%-15s", PQgetvalue(res, i, j));
    printf("\n");
}

PQclear(res);

/* close the portal ... we don't bother to check for errors ... */
res = PQexec(conn, "CLOSE myportal");
PQclear(res);

/* end the transaction */
res = PQexec(conn, "END");
PQclear(res);

/* close the connection to the database and cleanup */
PQfinish(conn);

return 0;
}
```

Example 35.2. libpq Example Program 2

```
/*
 * src/test/examples/testlibpq2.c
 *
 *
 * testlibpq2.c
 *     Test of the asynchronous notification interface
 *
 * Start this program, then from psql in another window do
 *     NOTIFY TBL2;
 * Repeat four times to get this program to exit.
 *
 * Or, if you want to get fancy, try this:
 * populate a database with the following commands
 * (provided in src/test/examples/testlibpq2.sql):
 *
 *     CREATE SCHEMA TESTLIBPQ2;
 *     SET search_path = TESTLIBPQ2;
 *     CREATE TABLE TBL1 (i int4);
 *     CREATE TABLE TBL2 (i int4);
 *     CREATE RULE r1 AS ON INSERT TO TBL1 DO
 *         (INSERT INTO TBL2 VALUES (new.i); NOTIFY TBL2);
 *
 * Start this program, then from psql do this four times:
 *
 *     INSERT INTO TESTLIBPQ2.TBL1 VALUES (10);
 */
```

```
#ifdef WIN32
#include <windows.h>
#endif
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <errno.h>
#include <sys/select.h>
#include <sys/time.h>
#include <sys/types.h>

#include "libpq-fe.h"

static void
exit_nicely(PGconn *conn)
{
    PQfinish(conn);
    exit(1);
}

int
main(int argc, char **argv)
{
    const char *conninfo;
    PGconn      *conn;
    PGresult     *res;
    PGnotify     *notify;
    int          nnotifies;

    /*
     * If the user supplies a parameter on the command line, use it as the
     * conninfo string; otherwise default to setting dbname=postgres and using
     * environment variables or defaults for all other connection parameters.
     */
    if (argc > 1)
        conninfo = argv[1];
    else
        conninfo = "dbname = postgres";

    /* Make a connection to the database */
    conn = PQconnectdb(conninfo);

    /* Check to see that the backend connection was successfully made */
    if (PQstatus(conn) != CONNECTION_OK)
    {
        fprintf(stderr, "%s", PQerrorMessage(conn));
        exit_nicely(conn);
    }

    /* Set always-secure search path, so malicious users can't take control. */
    res = PQexec(conn,
                  "SELECT pg_catalog.set_config('search_path', '', false)");
    if (PQresultStatus(res) != PGRES_TUPLES_OK)
    {
        fprintf(stderr, "SET failed: %s", PQerrorMessage(conn));
        PQclear(res);
        exit_nicely(conn);
    }
}
```

```
}

/*
 * Should PQclear PGresult whenever it is no longer needed to avoid memory
 * leaks
 */
PQclear(res);

/*
 * Issue LISTEN command to enable notifications from the rule's NOTIFY.
 */
res = PQexec(conn, "LISTEN TBL2");
if (PQresultStatus(res) != PGRES_COMMAND_OK)
{
    fprintf(stderr, "LISTEN command failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}
PQclear(res);

/* Quit after four notifies are received. */
nnotifies = 0;
while (nnotifies < 4)
{
    /*
     * Sleep until something happens on the connection. We use select(2)
     * to wait for input, but you could also use poll() or similar
     * facilities.
     */
    int          sock;
    fd_set       input_mask;

    sock = PQsocket(conn);

    if (sock < 0)
        break;                /* shouldn't happen */

    FD_ZERO(&input_mask);
    FD_SET(sock, &input_mask);

    if (select(sock + 1, &input_mask, NULL, NULL, NULL) < 0)
    {
        fprintf(stderr, "select() failed: %s\n", strerror(errno));
        exit_nicely(conn);
    }

    /* Now check for input */
    PQconsumeInput(conn);
    while ((notify = PQnotifies(conn)) != NULL)
    {
        fprintf(stderr,
            "ASYNC NOTIFY of '%s' received from backend PID %d\n",
            notify->relname, notify->be_pid);
        PQfreemem(notify);
        nnotifies++;
        PQconsumeInput(conn);
    }
}
}
```

```
fprintf(stderr, "Done.\n");

/* close the connection to the database and cleanup */
PQfinish(conn);

return 0;
}
```

Example 35.3. libpq Example Program 3

```
/*
 * src/test/examples/testlibpq3.c
 *
 *
 * testlibpq3.c
 *     Test out-of-line parameters and binary I/O.
 *
 * Before running this, populate a database with the following commands
 * (provided in src/test/examples/testlibpq3.sql):
 *
 * CREATE SCHEMA testlibpq3;
 * SET search_path = testlibpq3;
 * SET standard_conforming_strings = ON;
 * CREATE TABLE test1 (i int4, t text, b bytea);
 * INSERT INTO test1 values (1, 'joe's place', '\000\001\002\003\004');
 * INSERT INTO test1 values (2, 'ho there', '\004\003\002\001\000');
 *
 * The expected output is:
 *
 * tuple 0: got
 *   i = (4 bytes) 1
 *   t = (11 bytes) 'joe's place'
 *   b = (5 bytes) \000\001\002\003\004
 *
 * tuple 0: got
 *   i = (4 bytes) 2
 *   t = (8 bytes) 'ho there'
 *   b = (5 bytes) \004\003\002\001\000
 */

#ifdef WIN32
#include <windows.h>
#endif

#include <stdio.h>
#include <stdlib.h>
#include <stdint.h>
#include <string.h>
#include <sys/types.h>
#include "libpq-fe.h"

/* for ntohs/htons */
#include <netinet/in.h>
#include <arpa/inet.h>
```



```
static void
exit_nicely(PGconn *conn)
{
    PQfinish(conn);
    exit(1);
}

/*
 * This function prints a query result that is a binary-format fetch from
 * a table defined as in the comment above. We split it out because the
 * main() function uses it twice.
 */
static void
show_binary_results(PGresult *res)
{
    int        i,
               j;
    int        i_fnum,
               t_fnum,
               b_fnum;

    /* Use PQfnumber to avoid assumptions about field order in result */
    i_fnum = PQfnumber(res, "i");
    t_fnum = PQfnumber(res, "t");
    b_fnum = PQfnumber(res, "b");

    for (i = 0; i < PQntuples(res); i++)
    {
        char    *iptr;
        char    *tptr;
        char    *bptr;
        int      blen;
        int      ival;

        /* Get the field values (we ignore possibility they are null!) */
        iptr = PQgetvalue(res, i, i_fnum);
        tptr = PQgetvalue(res, i, t_fnum);
        bptr = PQgetvalue(res, i, b_fnum);

        /*
         * The binary representation of INT4 is in network byte order, which
         * we'd better coerce to the local byte order.
         */
        ival = ntohl(*(uint32_t *) iptr);

        /*
         * The binary representation of TEXT is, well, text, and since libpq
         * was nice enough to append a zero byte to it, it'll work just fine
         * as a C string.
         *
         * The binary representation of BYTEA is a bunch of bytes, which could
         * include embedded nulls so we have to pay attention to field length.
         */
        blen = PQgetlength(res, i, b_fnum);

        printf("tuple %d: got\n", i);
        printf(" i = (%d bytes) %d\n",
               PQgetlength(res, i, i_fnum), ival);
    }
}
```

```
        printf(" t = (%d bytes) '%s'\n",
               PQgetlength(res, i, t_fnum), tptr);
        printf(" b = (%d bytes) ", blen);
        for (j = 0; j < blen; j++)
            printf("\\%03o", bptr[j]);
        printf("\n\n");
    }
}

int
main(int argc, char **argv)
{
    const char *conninfo;
    PGconn      *conn;
    PGresult     *res;
    const char *paramValues[1];
    int          paramLengths[1];
    int          paramFormats[1];
    uint32_t     binaryIntVal;

    /*
     * If the user supplies a parameter on the command line, use it as the
     * conninfo string; otherwise default to setting dbname=postgres and using
     * environment variables or defaults for all other connection parameters.
     */
    if (argc > 1)
        conninfo = argv[1];
    else
        conninfo = "dbname = postgres";

    /* Make a connection to the database */
    conn = PQconnectdb(conninfo);

    /* Check to see that the backend connection was successfully made */
    if (PQstatus(conn) != CONNECTION_OK)
    {
        fprintf(stderr, "%s", PQerrorMessage(conn));
        exit_nicely(conn);
    }

    /* Set always-secure search path, so malicious users can't take control. */
    res = PQexec(conn, "SET search_path = testlibpq3");
    if (PQresultStatus(res) != PGRES_COMMAND_OK)
    {
        fprintf(stderr, "SET failed: %s", PQerrorMessage(conn));
        PQclear(res);
        exit_nicely(conn);
    }
    PQclear(res);

    /*
     * The point of this program is to illustrate use of PQexecParams() with
     * out-of-line parameters, as well as binary transmission of data.
     *
     * This first example transmits the parameters as text, but receives the
     * results in binary format. By using out-of-line parameters we can avoid
     * a lot of tedious mucking about with quoting and escaping, even though
     * the data is text. Notice how we don't have to do anything special with

```

```
* the quote mark in the parameter value.
*/

/* Here is our out-of-line parameter value */
paramValues[0] = "joe's place";

res = PQexecParams(conn,
    "SELECT * FROM test1 WHERE t = $1",
    1,          /* one param */
    NULL,      /* let the backend deduce param type */
    paramValues,
    NULL,      /* don't need param lengths since text */
    NULL,      /* default to all text params */
    1);        /* ask for binary results */

if (PQresultStatus(res) != PGRES_TUPLES_OK)
{
    fprintf(stderr, "SELECT failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}

show_binary_results(res);

PQclear(res);

/*
 * In this second example we transmit an integer parameter in binary form,
 * and again retrieve the results in binary form.
 *
 * Although we tell PQexecParams we are letting the backend deduce
 * parameter type, we really force the decision by casting the parameter
 * symbol in the query text. This is a good safety measure when sending
 * binary parameters.
 */

/* Convert integer value "2" to network byte order */
binaryIntVal = htonl((uint32_t) 2);

/* Set up parameter arrays for PQexecParams */
paramValues[0] = (char *) &binaryIntVal;
paramLengths[0] = sizeof(binaryIntVal);
paramFormats[0] = 1;          /* binary */

res = PQexecParams(conn,
    "SELECT * FROM test1 WHERE i = $1::int4",
    1,          /* one param */
    NULL,      /* let the backend deduce param type */
    paramValues,
    paramLengths,
    paramFormats,
    1);        /* ask for binary results */

if (PQresultStatus(res) != PGRES_TUPLES_OK)
{
    fprintf(stderr, "SELECT failed: %s", PQerrorMessage(conn));
    PQclear(res);
    exit_nicely(conn);
}
```

```
    }

    show_binary_results(res);

    PQclear(res);

    /* close the connection to the database and cleanup */
    PQfinish(conn);

    return 0;
}
```

Chapter 36. Large Objects

Postgres Pro has a *large object* facility, which provides stream-style access to user data that is stored in a special large-object structure. Streaming access is useful when working with data values that are too large to manipulate conveniently as a whole.

This chapter describes the implementation and the programming and query language interfaces to Postgres Pro large object data. We use the libpq C library for the examples in this chapter, but most programming interfaces native to Postgres Pro support equivalent functionality. Other interfaces might use the large object interface internally to provide generic support for large values. This is not described here.

36.1. Introduction

All large objects are stored in a single system table named `pg_largeobject`. Each large object also has an entry in the system table `pg_largeobject_metadata`. Large objects can be created, modified, and deleted using a read/write API that is similar to standard operations on files.

Postgres Pro also supports a storage system called “**TOAST**”, which automatically stores values larger than a single database page into a secondary storage area per table. This makes the large object facility partially obsolete. One remaining advantage of the large object facility is that it allows values up to 4 TB in size, whereas TOASTed fields can be at most 1 GB. Also, reading and updating portions of a large object can be done efficiently, while most operations on a TOASTed field will read or write the whole value as a unit.

36.2. Implementation Features

The large object implementation breaks large objects up into “chunks” and stores the chunks in rows in the database. A B-tree index guarantees fast searches for the correct chunk number when doing random access reads and writes.

The chunks stored for a large object do not have to be contiguous. For example, if an application opens a new large object, seeks to offset 1000000, and writes a few bytes there, this does not result in allocation of 1000000 bytes worth of storage; only of chunks covering the range of data bytes actually written. A read operation will, however, read out zeroes for any unallocated locations preceding the last existing chunk. This corresponds to the common behavior of “sparsely allocated” files in Unix file systems.

As of PostgreSQL 9.0, large objects have an owner and a set of access permissions, which can be managed using [GRANT](#) and [REVOKE](#). `SELECT` privileges are required to read a large object, and `UPDATE` privileges are required to write or truncate it. Only the large object's owner (or a database superuser) can delete, comment on, or change the owner of a large object. To adjust this behavior for compatibility with prior releases, see the [lo_compat_privileges](#) run-time parameter.

36.3. Client Interfaces

This section describes the facilities that Postgres Pro's libpq client interface library provides for accessing large objects. The Postgres Pro large object interface is modeled after the Unix file-system interface, with analogues of `open`, `read`, `write`, `lseek`, etc.

All large object manipulation using these functions *must* take place within an SQL transaction block, since large object file descriptors are only valid for the duration of a transaction. Write operations, including `lo_open` with the `INV_WRITE` mode, are not allowed in a read-only transaction.

If an error occurs while executing any one of these functions, the function will return an otherwise-impossible value, typically 0 or -1. A message describing the error is stored in the connection object and can be retrieved with [PQerrorMessage](#).

Client applications that use these functions should include the header file `libpq/libpq-fs.h` and link with the libpq library.

Client applications cannot use these functions while a libpq connection is in pipeline mode.

36.3.1. Creating a Large Object

The function

```
Oid lo_create(PGconn *conn, Oid lobjId);
```

creates a new large object. The OID to be assigned can be specified by *lobjId*; if so, failure occurs if that OID is already in use for some large object. If *lobjId* is `InvalidOid` (zero) then `lo_create` assigns an unused OID. The return value is the OID that was assigned to the new large object, or `InvalidOid` (zero) on failure.

An example:

```
inv_oid = lo_create(conn, desired_oid);
```

The older function

```
Oid lo_creat(PGconn *conn, int mode);
```

also creates a new large object, always assigning an unused OID. The return value is the OID that was assigned to the new large object, or `InvalidOid` (zero) on failure.

In PostgreSQL releases 8.1 and later, the *mode* is ignored, so that `lo_creat` is exactly equivalent to `lo_create` with a zero second argument. However, there is little reason to use `lo_creat` unless you need to work with servers older than 8.1. To work with such an old server, you must use `lo_creat` not `lo_create`, and you must set *mode* to one of `INV_READ`, `INV_WRITE`, or `INV_READ | INV_WRITE`. (These symbolic constants are defined in the header file `libpq/libpq-fs.h`.)

An example:

```
inv_oid = lo_creat(conn, INV_READ|INV_WRITE);
```

36.3.2. Importing a Large Object

To import an operating system file as a large object, call

```
Oid lo_import(PGconn *conn, const char *filename);
```

filename specifies the operating system name of the file to be imported as a large object. The return value is the OID that was assigned to the new large object, or `InvalidOid` (zero) on failure. Note that the file is read by the client interface library, not by the server; so it must exist in the client file system and be readable by the client application.

The function

```
Oid lo_import_with_oid(PGconn *conn, const char *filename, Oid lobjId);
```

also imports a new large object. The OID to be assigned can be specified by *lobjId*; if so, failure occurs if that OID is already in use for some large object. If *lobjId* is `InvalidOid` (zero) then `lo_import_with_oid` assigns an unused OID (this is the same behavior as `lo_import`). The return value is the OID that was assigned to the new large object, or `InvalidOid` (zero) on failure.

`lo_import_with_oid` is new as of PostgreSQL 8.4 and uses `lo_create` internally which is new in 8.1; if this function is run against 8.0 or before, it will fail and return `InvalidOid`.

36.3.3. Exporting a Large Object

To export a large object into an operating system file, call

```
int lo_export(PGconn *conn, Oid lobjId, const char *filename);
```

The *lobjId* argument specifies the OID of the large object to export and the *filename* argument specifies the operating system name of the file. Note that the file is written by the client interface library, not by the server. Returns 1 on success, -1 on failure.

36.3.4. Opening an Existing Large Object

To open an existing large object for reading or writing, call

```
int lo_open(PGconn *conn, Oid lobjId, int mode);
```

The `lobjId` argument specifies the OID of the large object to open. The `mode` bits control whether the object is opened for reading (`INV_READ`), writing (`INV_WRITE`), or both. (These symbolic constants are defined in the header file `libpq/libpq-fs.h`.) `lo_open` returns a (non-negative) large object descriptor for later use in `lo_read`, `lo_write`, `lo_lseek`, `lo_lseek64`, `lo_tell`, `lo_tell64`, `lo_truncate`, `lo_truncate64`, and `lo_close`. The descriptor is only valid for the duration of the current transaction. On failure, -1 is returned.

The server currently does not distinguish between modes `INV_WRITE` and `INV_READ | INV_WRITE`: you are allowed to read from the descriptor in either case. However there is a significant difference between these modes and `INV_READ` alone: with `INV_READ` you cannot write on the descriptor, and the data read from it will reflect the contents of the large object at the time of the transaction snapshot that was active when `lo_open` was executed, regardless of later writes by this or other transactions. Reading from a descriptor opened with `INV_WRITE` returns data that reflects all writes of other committed transactions as well as writes of the current transaction. This is similar to the behavior of `REPEATABLE READ` versus `READ COMMITTED` transaction modes for ordinary SQL `SELECT` commands.

`lo_open` will fail if `SELECT` privilege is not available for the large object, or if `INV_WRITE` is specified and `UPDATE` privilege is not available. (Prior to PostgreSQL 11, these privilege checks were instead performed at the first actual read or write call using the descriptor.) These privilege checks can be disabled with the [lo_compat_privileges](#) run-time parameter.

An example:

```
inv_fd = lo_open(conn, inv_oid, INV_READ|INV_WRITE);
```

36.3.5. Writing Data to a Large Object

The function

```
int lo_write(PGconn *conn, int fd, const char *buf, size_t len);
```

writes `len` bytes from `buf` (which must be of size `len`) to large object descriptor `fd`. The `fd` argument must have been returned by a previous `lo_open`. The number of bytes actually written is returned (in the current implementation, this will always equal `len` unless there is an error). In the event of an error, the return value is -1.

Although the `len` parameter is declared as `size_t`, this function will reject length values larger than `INT_MAX`. In practice, it's best to transfer data in chunks of at most a few megabytes anyway.

36.3.6. Reading Data from a Large Object

The function

```
int lo_read(PGconn *conn, int fd, char *buf, size_t len);
```

reads up to `len` bytes from large object descriptor `fd` into `buf` (which must be of size `len`). The `fd` argument must have been returned by a previous `lo_open`. The number of bytes actually read is returned; this will be less than `len` if the end of the large object is reached first. In the event of an error, the return value is -1.

Although the `len` parameter is declared as `size_t`, this function will reject length values larger than `INT_MAX`. In practice, it's best to transfer data in chunks of at most a few megabytes anyway.

36.3.7. Seeking in a Large Object

To change the current read or write location associated with a large object descriptor, call

```
int lo_lseek(PGconn *conn, int fd, int offset, int whence);
```

This function moves the current location pointer for the large object descriptor identified by `fd` to the new location specified by `offset`. The valid values for `whence` are `SEEK_SET` (seek from object start), `SEEK_CUR` (seek from current position), and `SEEK_END` (seek from object end). The return value is the new location pointer, or -1 on error.

When dealing with large objects that might exceed 2GB in size, instead use

```
pg_int64 lo_lseek64(PGconn *conn, int fd, pg_int64 offset, int whence);
```

This function has the same behavior as `lo_lseek`, but it can accept an *offset* larger than 2GB and/or deliver a result larger than 2GB. Note that `lo_lseek` will fail if the new location pointer would be greater than 2GB.

`lo_lseek64` is new as of PostgreSQL 9.3. If this function is run against an older server version, it will fail and return -1.

36.3.8. Obtaining the Seek Position of a Large Object

To obtain the current read or write location of a large object descriptor, call

```
int lo_tell(PGconn *conn, int fd);
```

If there is an error, the return value is -1.

When dealing with large objects that might exceed 2GB in size, instead use

```
pg_int64 lo_tell64(PGconn *conn, int fd);
```

This function has the same behavior as `lo_tell`, but it can deliver a result larger than 2GB. Note that `lo_tell` will fail if the current read/write location is greater than 2GB.

`lo_tell64` is new as of PostgreSQL 9.3. If this function is run against an older server version, it will fail and return -1.

36.3.9. Truncating a Large Object

To truncate a large object to a given length, call

```
int lo_truncate(PGconn *conn, int fd, size_t len);
```

This function truncates the large object descriptor *fd* to length *len*. The *fd* argument must have been returned by a previous `lo_open`. If *len* is greater than the large object's current length, the large object is extended to the specified length with null bytes ('\0'). On success, `lo_truncate` returns zero. On error, the return value is -1.

The read/write location associated with the descriptor *fd* is not changed.

Although the *len* parameter is declared as `size_t`, `lo_truncate` will reject length values larger than `INT_MAX`.

When dealing with large objects that might exceed 2GB in size, instead use

```
int lo_truncate64(PGconn *conn, int fd, pg_int64 len);
```

This function has the same behavior as `lo_truncate`, but it can accept a *len* value exceeding 2GB.

`lo_truncate` is new as of PostgreSQL 8.3; if this function is run against an older server version, it will fail and return -1.

`lo_truncate64` is new as of PostgreSQL 9.3; if this function is run against an older server version, it will fail and return -1.

36.3.10. Closing a Large Object Descriptor

A large object descriptor can be closed by calling

```
int lo_close(PGconn *conn, int fd);
```

where *fd* is a large object descriptor returned by `lo_open`. On success, `lo_close` returns zero. On error, the return value is -1.

Any large object descriptors that remain open at the end of a transaction will be closed automatically.

36.3.11. Removing a Large Object

To remove a large object from the database, call


```
int lo_unlink(PGconn *conn, Oid lobjId);
```

The *lobjId* argument specifies the OID of the large object to remove. Returns 1 if successful, -1 on failure.

36.4. Server-Side Functions

Server-side functions tailored for manipulating large objects from SQL are listed in [Table 36.1](#).

Table 36.1. SQL-Oriented Large Object Functions

Function	Description	Example(s)
<code>lo_from_bytea (loid oid, data bytea) → oid</code>	Creates a large object and stores <i>data</i> in it. If <i>loid</i> is zero then the system will choose a free OID, otherwise that OID is used (with an error if some large object already has that OID). On success, the large object's OID is returned.	<code>lo_from_bytea(0, '\xffffffff00') → 24528</code>
<code>lo_put (loid oid, offset bigint, data bytea) → void</code>	Writes <i>data</i> starting at the given offset within the large object; the large object is enlarged if necessary.	<code>lo_put(24528, 1, '\xaa') →</code>
<code>lo_get (loid oid [, offset bigint, length integer]) → bytea</code>	Extracts the large object's contents, or a substring thereof.	<code>lo_get(24528, 0, 3) → \xffaaff</code>

There are additional server-side functions corresponding to each of the client-side functions described earlier; indeed, for the most part the client-side functions are simply interfaces to the equivalent server-side functions. The ones just as convenient to call via SQL commands are `lo_creat`, `lo_create`, `lo_unlink`, `lo_import`, and `lo_export`. Here are examples of their use:

```
CREATE TABLE image (
    name          text,
    raster        oid
);

SELECT lo_creat(-1);           -- returns OID of new, empty large object

SELECT lo_create(43213);      -- attempts to create large object with OID 43213

SELECT lo_unlink(173454);     -- deletes large object with OID 173454

INSERT INTO image (name, raster)
VALUES ('beautiful image', lo_import('/etc/motd'));

INSERT INTO image (name, raster) -- same as above, but specify OID to use
VALUES ('beautiful image', lo_import('/etc/motd', 68583));

SELECT lo_export(image.raster, '/tmp/motd') FROM image
WHERE name = 'beautiful image';
```

The server-side `lo_import` and `lo_export` functions behave considerably differently from their client-side analogs. These two functions read and write files in the server's file system, using the permissions of the database's owning user. Therefore, by default their use is restricted to superusers. In contrast, the client-side import and export functions read and write files in the client's file system, using the permissions of the client program. The client-side functions do not require any database privileges, except the privilege to read or write the large object in question.

Caution

It is possible to **GRANT** use of the server-side `lo_import` and `lo_export` functions to non-superusers, but careful consideration of the security implications is required. A malicious user of such privileges could easily parlay them into becoming superuser (for example by rewriting server configuration files), or could attack the rest of the server's file system without bothering to obtain database superuser privileges as such. *Access to roles having such privilege must therefore be guarded just as carefully as access to superuser roles.* Nonetheless, if use of server-side `lo_import` or `lo_export` is needed for some routine task, it's safer to use a role with such privileges than one with full superuser privileges, as that helps to reduce the risk of damage from accidental errors.

The functionality of `lo_read` and `lo_write` is also available via server-side calls, but the names of the server-side functions differ from the client side interfaces in that they do not contain underscores. You must call these functions as `loread` and `lowrite`.

36.5. Example Program

[Example 36.1](#) is a sample program which shows how the large object interface in `libpq` can be used. Parts of the program are commented out for the reader's benefit.

Example 36.1. Large Objects with `libpq` Example Program

```
/*-----
 *
 * testlo.c
 *   test using large objects with libpq
 *
 * Portions Copyright (c) 1996-2024, PostgreSQL Global Development Group
 * Portions Copyright (c) 1994, Regents of the University of California
 *
 * IDENTIFICATION
 *   src/test/examples/testlo.c
 *-----
 */
#include <stdio.h>
#include <stdlib.h>

#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <unistd.h>

#include "libpq-fe.h"
#include "libpq/libpq-fs.h"

#define BUFSIZE          1024

/*
 * importFile -
 *   import file "in_filename" into database as large object "lobjOid"
 *
 */
static Oid
importFile(PGconn *conn, char *filename)
{
    Oid          lobjId;
```

```
int      lobj_fd;
char     buf[BUFSIZE];
int      nbytes,
        tmp;
int      fd;

/*
 * open the file to be read in
 */
fd = open(filename, O_RDONLY, 0666);
if (fd < 0)
{
    /* error */
    fprintf(stderr, "cannot open unix file\"%s\"\\n", filename);
}

/*
 * create the large object
 */
lobjId = lo_creat(conn, INV_READ | INV_WRITE);
if (lobjId == 0)
    fprintf(stderr, "cannot create large object");

lobj_fd = lo_open(conn, lobjId, INV_WRITE);

/*
 * read in from the Unix file and write to the inversion file
 */
while ((nbytes = read(fd, buf, BUFSIZE)) > 0)
{
    tmp = lo_write(conn, lobj_fd, buf, nbytes);
    if (tmp < nbytes)
        fprintf(stderr, "error while reading \"%s\"\\n", filename);
}

close(fd);
lo_close(conn, lobj_fd);

return lobjId;
}

static void
pickout(PGconn *conn, Oid lobjId, int start, int len)
{
    int      lobj_fd;
    char     *buf;
    int      nbytes;
    int      nread;

    lobj_fd = lo_open(conn, lobjId, INV_READ);
    if (lobj_fd < 0)
        fprintf(stderr, "cannot open large object %u", lobjId);

    lo_lseek(conn, lobj_fd, start, SEEK_SET);
    buf = malloc(len + 1);

    nread = 0;
    while (len - nread > 0)
    {
```

```
        nbytes = lo_read(conn, lobj_fd, buf, len - nread);
        buf[nbytes] = '\0';
        fprintf(stderr, ">>> %s", buf);
        nread += nbytes;
        if (nbytes <= 0)
            break;                /* no more data? */
    }
    free(buf);
    fprintf(stderr, "\n");
    lo_close(conn, lobj_fd);
}

static void
overwrite(PGconn *conn, Oid lobjId, int start, int len)
{
    int         lobj_fd;
    char        *buf;
    int         nbytes;
    int         nwritten;
    int         i;

    lobj_fd = lo_open(conn, lobjId, INV_WRITE);
    if (lobj_fd < 0)
        fprintf(stderr, "cannot open large object %u", lobjId);

    lo_lseek(conn, lobj_fd, start, SEEK_SET);
    buf = malloc(len + 1);

    for (i = 0; i < len; i++)
        buf[i] = 'X';
    buf[i] = '\0';

    nwritten = 0;
    while (len - nwritten > 0)
    {
        nbytes = lo_write(conn, lobj_fd, buf + nwritten, len - nwritten);
        nwritten += nbytes;
        if (nbytes <= 0)
        {
            fprintf(stderr, "\nWRITE FAILED!\n");
            break;
        }
    }
    free(buf);
    fprintf(stderr, "\n");
    lo_close(conn, lobj_fd);
}

/*
 * exportFile -
 *     export large object "lobjOid" to file "out_filename"
 */
static void
exportFile(PGconn *conn, Oid lobjId, char *filename)
{
    int         lobj_fd;
```

```
char        buf[BUFSIZE];
int         nbytes,
           tmp;
int         fd;

/*
 * open the large object
 */
lobj_fd = lo_open(conn, lobjId, INV_READ);
if (lobj_fd < 0)
    fprintf(stderr, "cannot open large object %u", lobjId);

/*
 * open the file to be written to
 */
fd = open(filename, O_CREAT | O_WRONLY | O_TRUNC, 0666);
if (fd < 0)
{
    /* error */
    fprintf(stderr, "cannot open unix file\"%s\"",
            filename);
}

/*
 * read in from the inversion file and write to the Unix file
 */
while ((nbytes = lo_read(conn, lobj_fd, buf, BUFSIZE)) > 0)
{
    tmp = write(fd, buf, nbytes);
    if (tmp < nbytes)
    {
        fprintf(stderr, "error while writing \"%s\"",
                filename);
    }
}

lo_close(conn, lobj_fd);
close(fd);
}

static void
exit_nicely(PGconn *conn)
{
    PQfinish(conn);
    exit(1);
}

int
main(int argc, char **argv)
{
    char        *in_filename,
               *out_filename;
    char        *database;
    Oid         lobjOid;
    PGconn      *conn;
    PGresult     *res;

    if (argc != 4)
    {
```

```
        fprintf(stderr, "Usage: %s database_name in_filename out_filename\n",
                argv[0]);
        exit(1);
    }

    database = argv[1];
    in_filename = argv[2];
    out_filename = argv[3];

    /*
     * set up the connection
     */
    conn = PQsetdb(NULL, NULL, NULL, NULL, database);

    /* check to see that the backend connection was successfully made */
    if (PQstatus(conn) != CONNECTION_OK)
    {
        fprintf(stderr, "%s", PQerrorMessage(conn));
        exit_nicely(conn);
    }

    /* Set always-secure search path, so malicious users can't take control. */
    res = PQexec(conn,
                  "SELECT pg_catalog.set_config('search_path', '', false)");
    if (PQresultStatus(res) != PGRES_TUPLES_OK)
    {
        fprintf(stderr, "SET failed: %s", PQerrorMessage(conn));
        PQclear(res);
        exit_nicely(conn);
    }
    PQclear(res);

    res = PQexec(conn, "begin");
    PQclear(res);
    printf("importing file \"%s\" ...\n", in_filename);
    /* lobjOid = importFile(conn, in_filename); */
    lobjOid = lo_import(conn, in_filename);
    if (lobjOid == 0)
        fprintf(stderr, "%s\n", PQerrorMessage(conn));
    else
    {
        printf("\tas large object %u.\n", lobjOid);

        printf("picking out bytes 1000-2000 of the large object\n");
        pickout(conn, lobjOid, 1000, 1000);

        printf("overwriting bytes 1000-2000 of the large object with X's\n");
        overwrite(conn, lobjOid, 1000, 1000);

        printf("exporting large object to file \"%s\" ...\n", out_filename);
        /* exportFile(conn, lobjOid, out_filename); */
        if (lo_export(conn, lobjOid, out_filename) < 0)
            fprintf(stderr, "%s\n", PQerrorMessage(conn));
    }

    res = PQexec(conn, "end");
    PQclear(res);
    PQfinish(conn);
```

```
    return 0;  
}
```

Chapter 37. ECPG — Embedded SQL in C

This chapter describes the embedded SQL package for Postgres Pro. It was written by Linus Tolke (<linus@epact.se>) and Michael Meskes (<meskes@postgresql.org>). Originally it was written to work with C. It also works with C++, but it does not recognize all C++ constructs yet.

This documentation is quite incomplete. But since this interface is standardized, additional information can be found in many resources about SQL.

37.1. The Concept

An embedded SQL program consists of code written in an ordinary programming language, in this case C, mixed with SQL commands in specially marked sections. To build the program, the source code (*.pgc) is first passed through the embedded SQL preprocessor, which converts it to an ordinary C program (*.c), and afterwards it can be processed by a C compiler. (For details about the compiling and linking see [Section 37.10](#).) Converted ECPG applications call functions in the libpq library through the embedded SQL library (ecpglib), and communicate with the Postgres Pro server using the normal frontend-backend protocol.

Embedded SQL has advantages over other methods for handling SQL commands from C code. First, it takes care of the tedious passing of information to and from variables in your C program. Second, the SQL code in the program is checked at build time for syntactical correctness. Third, embedded SQL in C is specified in the SQL standard and supported by many other SQL database systems. The Postgres Pro implementation is designed to match this standard as much as possible, and it is usually possible to port embedded SQL programs written for other SQL databases to Postgres Pro with relative ease.

As already stated, programs written for the embedded SQL interface are normal C programs with special code inserted to perform database-related actions. This special code always has the form:

```
EXEC SQL ...;
```

These statements syntactically take the place of a C statement. Depending on the particular statement, they can appear at the global level or within a function.

Embedded SQL statements follow the case-sensitivity rules of normal SQL code, and not those of C. Also they allow nested C-style comments as per the SQL standard. The C part of the program, however, follows the C standard of not accepting nested comments. Embedded SQL statements likewise use SQL rules, not C rules, for parsing quoted strings and identifiers. (See [Section 4.1.2.1](#) and [Section 4.1.1](#) respectively. Note that ECPG assumes that `standard_conforming_strings` is on.) Of course, the C part of the program follows C quoting rules.

The following sections explain all the embedded SQL statements.

37.2. Managing Database Connections

This section describes how to open, close, and switch database connections.

37.2.1. Connecting to the Database Server

One connects to a database using the following statement:

```
EXEC SQL CONNECT TO target [AS connection-name] [USER user-name];
```

The *target* can be specified in the following ways:

- `dbname[@hostname][:port]`
- `tcp:postgresql://hostname[:port][/dbname][?options]`
- `unix:postgresql://localhost[:port][/dbname][?options]`
- an SQL string literal containing one of the above forms
- a reference to a character variable containing one of the above forms (see examples)

- `DEFAULT`

The connection target `DEFAULT` initiates a connection to the default database under the default user name. No separate user name or connection name can be specified in that case.

If you specify the connection target directly (that is, not as a string literal or variable reference), then the components of the target are passed through normal SQL parsing; this means that, for example, the *hostname* must look like one or more SQL identifiers separated by dots, and those identifiers will be case-folded unless double-quoted. Values of any *options* must be SQL identifiers, integers, or variable references. Of course, you can put nearly anything into an SQL identifier by double-quoting it. In practice, it is probably less error-prone to use a (single-quoted) string literal or a variable reference than to write the connection target directly.

There are also different ways to specify the user name:

- `username`
- `username/password`
- `username IDENTIFIED BY password`
- `username USING password`

As above, the parameters *username* and *password* can be an SQL identifier, an SQL string literal, or a reference to a character variable.

If the connection target includes any *options*, those consist of *keyword=value* specifications separated by ampersands (&). The allowed key words are the same ones recognized by libpq (see [Section 35.1.2](#)). Spaces are ignored before any *keyword* or *value*, though not within or after one. Note that there is no way to write & within a *value*.

Notice that when specifying a socket connection (with the `unix:` prefix), the host name must be exactly `localhost`. To select a non-default socket directory, write the directory's pathname as the value of a *host* option in the *options* part of the target.

The *connection-name* is used to handle multiple connections in one program. It can be omitted if a program uses only one connection. The most recently opened connection becomes the current connection, which is used by default when an SQL statement is to be executed (see later in this chapter).

Here are some examples of `CONNECT` statements:

```
EXEC SQL CONNECT TO mydb@sql.mydomain.com;

EXEC SQL CONNECT TO tcp:postgresql://sql.mydomain.com/mydb AS myconnection USER john;

EXEC SQL BEGIN DECLARE SECTION;
const char *target = "mydb@sql.mydomain.com";
const char *user = "john";
const char *passwd = "secret";
EXEC SQL END DECLARE SECTION;
...
EXEC SQL CONNECT TO :target USER :user USING :passwd;
/* or EXEC SQL CONNECT TO :target USER :user/:passwd; */
```

The last example makes use of the feature referred to above as character variable references. You will see in later sections how C variables can be used in SQL statements when you prefix them with a colon.

Be advised that the format of the connection target is not specified in the SQL standard. So if you want to develop portable applications, you might want to use something based on the last example above to encapsulate the connection target string somewhere.

If untrusted users have access to a database that has not adopted a [secure schema usage pattern](#), begin each session by removing publicly-writable schemas from `search_path`. For example, add `options=-c`

`search_path=` to *options*, or issue `EXEC SQL SELECT pg_catalog.set_config('search_path', '', false);` after connecting. This consideration is not specific to ECPG; it applies to every interface for executing arbitrary SQL commands.

37.2.2. Choosing a Connection

SQL statements in embedded SQL programs are by default executed on the current connection, that is, the most recently opened one. If an application needs to manage multiple connections, then there are three ways to handle this.

The first option is to explicitly choose a connection for each SQL statement, for example:

```
EXEC SQL AT connection-name SELECT ...;
```

This option is particularly suitable if the application needs to use several connections in mixed order.

If your application uses multiple threads of execution, they cannot share a connection concurrently. You must either explicitly control access to the connection (using mutexes) or use a connection for each thread.

The second option is to execute a statement to switch the current connection. That statement is:

```
EXEC SQL SET CONNECTION connection-name;
```

This option is particularly convenient if many statements are to be executed on the same connection.

Here is an example program managing multiple database connections:

```
#include <stdio.h>

EXEC SQL BEGIN DECLARE SECTION;
    char dbname[1024];
EXEC SQL END DECLARE SECTION;

int
main()
{
    EXEC SQL CONNECT TO testdb1 AS con1 USER testuser;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
    EXEC SQL CONNECT TO testdb2 AS con2 USER testuser;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
    EXEC SQL CONNECT TO testdb3 AS con3 USER testuser;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    /* This query would be executed in the last opened database "testdb3". */
    EXEC SQL SELECT current_database() INTO :dbname;
    printf("current=%s (should be testdb3)\n", dbname);

    /* Using "AT" to run a query in "testdb2" */
    EXEC SQL AT con2 SELECT current_database() INTO :dbname;
    printf("current=%s (should be testdb2)\n", dbname);

    /* Switch the current connection to "testdb1". */
    EXEC SQL SET CONNECTION con1;

    EXEC SQL SELECT current_database() INTO :dbname;
    printf("current=%s (should be testdb1)\n", dbname);

    EXEC SQL DISCONNECT ALL;
    return 0;
}
```

This example would produce this output:

```
current=testdb3 (should be testdb3)
current=testdb2 (should be testdb2)
current=testdb1 (should be testdb1)
```

The third option is to declare an SQL identifier linked to the connection, for example:

```
EXEC SQL AT connection-name DECLARE statement-name STATEMENT;
EXEC SQL PREPARE statement-name FROM :dyn-string;
```

Once you link an SQL identifier to a connection, you execute dynamic SQL without an AT clause. Note that this option behaves like preprocessor directives, therefore the link is enabled only in the file.

Here is an example program using this option:

```
#include <stdio.h>

EXEC SQL BEGIN DECLARE SECTION;
char dbname[128];
char *dyn_sql = "SELECT current_database()";
EXEC SQL END DECLARE SECTION;

int main(){
    EXEC SQL CONNECT TO postgres AS con1;
    EXEC SQL CONNECT TO testdb AS con2;
    EXEC SQL AT con1 DECLARE stmt STATEMENT;
    EXEC SQL PREPARE stmt FROM :dyn_sql;
    EXEC SQL EXECUTE stmt INTO :dbname;
    printf("%s\n", dbname);

    EXEC SQL DISCONNECT ALL;
    return 0;
}
```

This example would produce this output, even if the default connection is testdb:

```
postgres
```

37.2.3. Closing a Connection

To close a connection, use the following statement:

```
EXEC SQL DISCONNECT [connection];
```

The *connection* can be specified in the following ways:

- *connection-name*
- CURRENT
- ALL

If no connection name is specified, the current connection is closed.

It is good style that an application always explicitly disconnect from every connection it opened.

37.3. Running SQL Commands

Any SQL command can be run from within an embedded SQL application. Below are some examples of how to do that.

37.3.1. Executing SQL Statements

Creating a table:

```
EXEC SQL CREATE TABLE foo (number integer, ascii char(16));
EXEC SQL CREATE UNIQUE INDEX num1 ON foo(number);
EXEC SQL COMMIT;
```

Inserting rows:

```
EXEC SQL INSERT INTO foo (number, ascii) VALUES (9999, 'doodad');
EXEC SQL COMMIT;
```

Deleting rows:

```
EXEC SQL DELETE FROM foo WHERE number = 9999;
EXEC SQL COMMIT;
```

Updates:

```
EXEC SQL UPDATE foo
    SET ascii = 'foobar'
    WHERE number = 9999;
EXEC SQL COMMIT;
```

`SELECT` statements that return a single result row can also be executed using `EXEC SQL` directly. To handle result sets with multiple rows, an application has to use a cursor; see [Section 37.3.2](#) below. (As a special case, an application can fetch multiple rows at once into an array host variable; see [Section 37.4.4.3.1](#).)

Single-row select:

```
EXEC SQL SELECT foo INTO :FooBar FROM table1 WHERE ascii = 'doodad';
```

Also, a configuration parameter can be retrieved with the `SHOW` command:

```
EXEC SQL SHOW search_path INTO :var;
```

The tokens of the form `:something` are *host variables*, that is, they refer to variables in the C program. They are explained in [Section 37.4](#).

37.3.2. Using Cursors

To retrieve a result set holding multiple rows, an application has to declare a cursor and fetch each row from the cursor. The steps to use a cursor are the following: declare a cursor, open it, fetch a row from the cursor, repeat, and finally close it.

Select using cursors:

```
EXEC SQL DECLARE foo_bar CURSOR FOR
    SELECT number, ascii FROM foo
    ORDER BY ascii;
EXEC SQL OPEN foo_bar;
EXEC SQL FETCH foo_bar INTO :FooBar, DooDad;
...
EXEC SQL CLOSE foo_bar;
EXEC SQL COMMIT;
```

For more details about declaring a cursor, see [DECLARE](#); for more details about fetching rows from a cursor, see [FETCH](#).

Note

The ECPG `DECLARE` command does not actually cause a statement to be sent to the Postgres Pro backend. The cursor is opened in the backend (using the backend's `DECLARE` command) at the point when the `OPEN` command is executed.

37.3.3. Managing Transactions

In the default mode, statements are committed only when `EXEC SQL COMMIT` is issued. The embedded SQL interface also supports autocommit of transactions (similar to `psql`'s default behavior) via the `-t` command-line option to `ecpg` (see [ecpg](#)) or via the `EXEC SQL SET AUTOCOMMIT TO ON` statement. In autocommit mode, each command is automatically committed unless it is inside an explicit transaction block. This mode can be explicitly turned off using `EXEC SQL SET AUTOCOMMIT TO OFF`.

The following transaction management commands are available:

```
EXEC SQL COMMIT
```

Commit an in-progress transaction.

```
EXEC SQL ROLLBACK
```

Roll back an in-progress transaction.

```
EXEC SQL PREPARE TRANSACTION transaction_id
```

Prepare the current transaction for two-phase commit.

```
EXEC SQL COMMIT PREPARED transaction_id
```

Commit a transaction that is in prepared state.

```
EXEC SQL ROLLBACK PREPARED transaction_id
```

Roll back a transaction that is in prepared state.

```
EXEC SQL SET AUTOCOMMIT TO ON
```

Enable autocommit mode.

```
EXEC SQL SET AUTOCOMMIT TO OFF
```

Disable autocommit mode. This is the default.

37.3.4. Prepared Statements

When the values to be passed to an SQL statement are not known at compile time, or the same statement is going to be used many times, then prepared statements can be useful.

The statement is prepared using the command `PREPARE`. For the values that are not known yet, use the placeholder `"?"`:

```
EXEC SQL PREPARE stmt1 FROM "SELECT oid, datname FROM pg_database WHERE oid = ?";
```

If a statement returns a single row, the application can call `EXECUTE` after `PREPARE` to execute the statement, supplying the actual values for the placeholders with a `USING` clause:

```
EXEC SQL EXECUTE stmt1 INTO :dboid, :dbname USING 1;
```

If a statement returns multiple rows, the application can use a cursor declared based on the prepared statement. To bind input parameters, the cursor must be opened with a `USING` clause:

```
EXEC SQL PREPARE stmt1 FROM "SELECT oid, datname FROM pg_database WHERE oid > ?";  
EXEC SQL DECLARE foo_bar CURSOR FOR stmt1;
```

```
/* when end of result set reached, break out of while loop */  
EXEC SQL WHENEVER NOT FOUND DO BREAK;
```

```
EXEC SQL OPEN foo_bar USING 100;
```

```
...  
while (1)
```

```
{  
    EXEC SQL FETCH NEXT FROM foo_bar INTO :dboid, :dbname;  
    ...  
}  
EXEC SQL CLOSE foo_bar;
```

When you don't need the prepared statement anymore, you should deallocate it:

```
EXEC SQL DEALLOCATE PREPARE name;
```

For more details about `PREPARE`, see [PREPARE](#). Also see [Section 37.5](#) for more details about using placeholders and input parameters.

37.4. Using Host Variables

In [Section 37.3](#) you saw how you can execute SQL statements from an embedded SQL program. Some of those statements only used fixed values and did not provide a way to insert user-supplied values into statements or have the program process the values returned by the query. Those kinds of statements are not really useful in real applications. This section explains in detail how you can pass data between your C program and the embedded SQL statements using a simple mechanism called *host variables*. In an embedded SQL program we consider the SQL statements to be *guests* in the C program code which is the *host language*. Therefore the variables of the C program are called *host variables*.

Another way to exchange values between Postgres Pro backends and ECPG applications is the use of SQL descriptors, described in [Section 37.7](#).

37.4.1. Overview

Passing data between the C program and the SQL statements is particularly simple in embedded SQL. Instead of having the program paste the data into the statement, which entails various complications, such as properly quoting the value, you can simply write the name of a C variable into the SQL statement, prefixed by a colon. For example:

```
EXEC SQL INSERT INTO sometable VALUES (:v1, 'foo', :v2);
```

This statement refers to two C variables named `v1` and `v2` and also uses a regular SQL string literal, to illustrate that you are not restricted to use one kind of data or the other.

This style of inserting C variables in SQL statements works anywhere a value expression is expected in an SQL statement.

37.4.2. Declare Sections

To pass data from the program to the database, for example as parameters in a query, or to pass data from the database back to the program, the C variables that are intended to contain this data need to be declared in specially marked sections, so the embedded SQL preprocessor is made aware of them.

This section starts with:

```
EXEC SQL BEGIN DECLARE SECTION;
```

and ends with:

```
EXEC SQL END DECLARE SECTION;
```

Between those lines, there must be normal C variable declarations, such as:

```
int    x = 4;  
char   foo[16], bar[16];
```

As you can see, you can optionally assign an initial value to the variable. The variable's scope is determined by the location of its declaring section within the program. You can also declare variables with the following syntax which implicitly creates a declare section:

```
EXEC SQL int i = 4;
```

You can have as many declare sections in a program as you like.

The declarations are also echoed to the output file as normal C variables, so there's no need to declare them again. Variables that are not intended to be used in SQL commands can be declared normally outside these special sections.

The definition of a structure or union also must be listed inside a `DECLARE` section. Otherwise the pre-processor cannot handle these types since it does not know the definition.

37.4.3. Retrieving Query Results

Now you should be able to pass data generated by your program into an SQL command. But how do you retrieve the results of a query? For that purpose, embedded SQL provides special variants of the usual commands `SELECT` and `FETCH`. These commands have a special `INTO` clause that specifies which host variables the retrieved values are to be stored in. `SELECT` is used for a query that returns only single row, and `FETCH` is used for a query that returns multiple rows, using a cursor.

Here is an example:

```
/*
 * assume this table:
 * CREATE TABLE test1 (a int, b varchar(50));
 */
```

```
EXEC SQL BEGIN DECLARE SECTION;
int v1;
VARCHAR v2;
EXEC SQL END DECLARE SECTION;
```

...

```
EXEC SQL SELECT a, b INTO :v1, :v2 FROM test;
```

So the `INTO` clause appears between the select list and the `FROM` clause. The number of elements in the select list and the list after `INTO` (also called the target list) must be equal.

Here is an example using the command `FETCH`:

```
EXEC SQL BEGIN DECLARE SECTION;
int v1;
VARCHAR v2;
EXEC SQL END DECLARE SECTION;
```

...

```
EXEC SQL DECLARE foo CURSOR FOR SELECT a, b FROM test;
```

...

```
do
{
    ...
    EXEC SQL FETCH NEXT FROM foo INTO :v1, :v2;
    ...
} while (...);
```

Here the `INTO` clause appears after all the normal clauses.

37.4.4. Type Mapping

When ECPG applications exchange values between the Postgres Pro server and the C application, such as when retrieving query results from the server or executing SQL statements with input parameters,

the values need to be converted between Postgres Pro data types and host language variable types (C language data types, concretely). One of the main points of ECPG is that it takes care of this automatically in most cases.

In this respect, there are two kinds of data types: Some simple Postgres Pro data types, such as `integer` and `text`, can be read and written by the application directly. Other Postgres Pro data types, such as `timestamp` and `numeric` can only be accessed through special library functions; see [Section 37.4.4.2](#).

[Table 37.1](#) shows which Postgres Pro data types correspond to which C data types. When you wish to send or receive a value of a given Postgres Pro data type, you should declare a C variable of the corresponding C data type in the declare section.

Table 37.1. Mapping Between Postgres Pro Data Types and C Variable Types

Postgres Pro data type	Host variable type
<code>smallint</code>	<code>short</code>
<code>integer</code>	<code>int</code>
<code>bigint</code>	<code>long long int</code>
<code>decimal</code>	<code>decimal</code> ^a
<code>numeric</code>	<code>numeric</code> ^a
<code>real</code>	<code>float</code>
<code>double precision</code>	<code>double</code>
<code>smallserial</code>	<code>short</code>
<code>serial</code>	<code>int</code>
<code>bigserial</code>	<code>long long int</code>
<code>oid</code>	<code>unsigned int</code>
<code>character(n), varchar(n), text</code>	<code>char[n+1], VARCHAR[n+1]</code>
<code>name</code>	<code>char[NAMEDATALEN]</code>
<code>timestamp</code>	<code>timestamp</code> ^a
<code>interval</code>	<code>interval</code> ^a
<code>date</code>	<code>date</code> ^a
<code>boolean</code>	<code>bool</code> ^b
<code>bytea</code>	<code>char *, bytea[n]</code>

^aThis type can only be accessed through special library functions; see [Section 37.4.4.2](#).

^bdeclared in `ecpglib.h` if not native

37.4.4.1. Handling Character Strings

To handle SQL character string data types, such as `varchar` and `text`, there are two possible ways to declare the host variables.

One way is using `char[]`, an array of `char`, which is the most common way to handle character data in C.

```
EXEC SQL BEGIN DECLARE SECTION;
    char str[50];
EXEC SQL END DECLARE SECTION;
```

Note that you have to take care of the length yourself. If you use this host variable as the target variable of a query which returns a string with more than 49 characters, a buffer overflow occurs.

The other way is using the `VARCHAR` type, which is a special type provided by ECPG. The definition on an array of type `VARCHAR` is converted into a named `struct` for every variable. A declaration like:

```
VARCHAR var[180];
```


is converted into:

```
struct varchar_var { int len; char arr[180]; } var;
```

The member `arr` hosts the string including a terminating zero byte. Thus, to store a string in a `VARCHAR` host variable, the host variable has to be declared with the length including the zero byte terminator. The member `len` holds the length of the string stored in the `arr` without the terminating zero byte. When a host variable is used as input for a query, if `strlen(arr)` and `len` are different, the shorter one is used.

`VARCHAR` can be written in upper or lower case, but not in mixed case.

`char` and `VARCHAR` host variables can also hold values of other SQL types, which will be stored in their string forms.

37.4.4.2. Accessing Special Data Types

ECPG contains some special types that help you to interact easily with some special data types from the Postgres Pro server. In particular, it has implemented support for the `numeric`, `decimal`, `date`, `timestamp`, and `interval` types. These data types cannot usefully be mapped to primitive host variable types (such as `int`, `long long int`, or `char[]`), because they have a complex internal structure. Applications deal with these types by declaring host variables in special types and accessing them using functions in the `pgtypes` library. The `pgtypes` library, described in detail in [Section 37.6](#) contains basic functions to deal with those types, such that you do not need to send a query to the SQL server just for adding an interval to a time stamp for example.

The follow subsections describe these special data types. For more details about `pgtypes` library functions, see [Section 37.6](#).

37.4.4.2.1. timestamp, date

Here is a pattern for handling `timestamp` variables in the ECPG host application.

First, the program has to include the header file for the `timestamp` type:

```
#include <pgtypes_timestamp.h>
```

Next, declare a host variable as type `timestamp` in the declare section:

```
EXEC SQL BEGIN DECLARE SECTION;
timestamp ts;
EXEC SQL END DECLARE SECTION;
```

And after reading a value into the host variable, process it using `pgtypes` library functions. In following example, the `timestamp` value is converted into text (ASCII) form with the `PGTYPEStimestamp_to_asc()` function:

```
EXEC SQL SELECT now()::timestamp INTO :ts;

printf("ts = %s\n", PGTYPEStimestamp_to_asc(ts));
```

This example will show some result like following:

```
ts = 2010-06-27 18:03:56.949343
```

In addition, the `DATE` type can be handled in the same way. The program has to include `pgtypes_date.h`, declare a host variable as the date type and convert a `DATE` value into a text form using `PGTYPES-date_to_asc()` function. For more details about the `pgtypes` library functions, see [Section 37.6](#).

37.4.4.2.2. interval

The handling of the `interval` type is also similar to the `timestamp` and `date` types. It is required, however, to allocate memory for an `interval` type value explicitly. In other words, the memory space for the variable has to be allocated in the heap memory, not in the stack memory.

Here is an example program:

```
#include <stdio.h>
#include <stdlib.h>
#include <pgtypes_interval.h>

int
main(void)
{
EXEC SQL BEGIN DECLARE SECTION;
    interval *in;
EXEC SQL END DECLARE SECTION;

    EXEC SQL CONNECT TO testdb;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    in = PGTYPESEinterval_new();
    EXEC SQL SELECT '1 min'::interval INTO :in;
    printf("interval = %s\n", PGTYPESEinterval_to_asc(in));
    PGTYPESEinterval_free(in);

    EXEC SQL COMMIT;
    EXEC SQL DISCONNECT ALL;
    return 0;
}
```

37.4.4.2.3. numeric, decimal

The handling of the `numeric` and `decimal` types is similar to the `interval` type: It requires defining a pointer, allocating some memory space on the heap, and accessing the variable using the `pgtypes` library functions. For more details about the `pgtypes` library functions, see [Section 37.6](#).

No functions are provided specifically for the `decimal` type. An application has to convert it to a `numeric` variable using a `pgtypes` library function to do further processing.

Here is an example program handling `numeric` and `decimal` type variables.

```
#include <stdio.h>
#include <stdlib.h>
#include <pgtypes_numeric.h>

EXEC SQL WHENEVER SQLERROR STOP;

int
main(void)
{
EXEC SQL BEGIN DECLARE SECTION;
    numeric *num;
    numeric *num2;
    decimal *dec;
EXEC SQL END DECLARE SECTION;

    EXEC SQL CONNECT TO testdb;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    num = PGTYPESEnumeric_new();
    dec = PGTYPESEdecimal_new();

    EXEC SQL SELECT 12.345::numeric(4,2), 23.456::decimal(4,2) INTO :num, :dec;
```

```
printf("numeric = %s\n", PGTYPEStnumeric_to_asc(num, 0));
printf("numeric = %s\n", PGTYPEStnumeric_to_asc(num, 1));
printf("numeric = %s\n", PGTYPEStnumeric_to_asc(num, 2));

/* Convert decimal to numeric to show a decimal value. */
num2 = PGTYPEStnumeric_new();
PGTYPEStnumeric_from_decimal(dec, num2);

printf("decimal = %s\n", PGTYPEStnumeric_to_asc(num2, 0));
printf("decimal = %s\n", PGTYPEStnumeric_to_asc(num2, 1));
printf("decimal = %s\n", PGTYPEStnumeric_to_asc(num2, 2));

PGTYPEStnumeric_free(num2);
PGTYPEStdecimal_free(dec);
PGTYPEStnumeric_free(num);

EXEC SQL COMMIT;
EXEC SQL DISCONNECT ALL;
return 0;
}
```

37.4.4.2.4. bytea

The handling of the `bytea` type is similar to that of `VARCHAR`. The definition on an array of type `bytea` is converted into a named struct for every variable. A declaration like:

```
bytea var[180];
```

is converted into:

```
struct bytea_var { int len; char arr[180]; } var;
```

The member `arr` hosts binary format data. It can also handle `'\0'` as part of data, unlike `VARCHAR`. The data is converted from/to hex format and sent/received by `ecpglib`.

Note

`bytea` variable can be used only when `bytea_output` is set to `hex`.

37.4.4.3. Host Variables with Nonprimitive Types

As a host variable you can also use arrays, typedefs, structs, and pointers.

37.4.4.3.1. Arrays

There are two use cases for arrays as host variables. The first is a way to store some text string in `char[]` or `VARCHAR[]`, as explained in [Section 37.4.4.1](#). The second use case is to retrieve multiple rows from a query result without using a cursor. Without an array, to process a query result consisting of multiple rows, it is required to use a cursor and the `FETCH` command. But with array host variables, multiple rows can be received at once. The length of the array has to be defined to be able to accommodate all rows, otherwise a buffer overflow will likely occur.

Following example scans the `pg_database` system table and shows all OIDs and names of the available databases:

```
int
main(void)
{
EXEC SQL BEGIN DECLARE SECTION;
    int dbid[8];
    char dbname[8][16];
```

```
int i;
EXEC SQL END DECLARE SECTION;

memset(dbname, 0, sizeof(char)* 16 * 8);
memset(dbid, 0, sizeof(int) * 8);

EXEC SQL CONNECT TO testdb;
EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

/* Retrieve multiple rows into arrays at once. */
EXEC SQL SELECT oid,datname INTO :dbid, :dbname FROM pg_database;

for (i = 0; i < 8; i++)
    printf("oid=%d, dbname=%s\n", dbid[i], dbname[i]);

EXEC SQL COMMIT;
EXEC SQL DISCONNECT ALL;
return 0;
}
```

This example shows following result. (The exact values depend on local circumstances.)

```
oid=1, dbname=template1
oid=11510, dbname=template0
oid=11511, dbname=postgres
oid=313780, dbname=testdb
oid=0, dbname=
oid=0, dbname=
oid=0, dbname=
```

37.4.4.3.2. Structures

A structure whose member names match the column names of a query result, can be used to retrieve multiple columns at once. The structure enables handling multiple column values in a single host variable.

The following example retrieves OIDs, names, and sizes of the available databases from the `pg_database` system table and using the `pg_database_size()` function. In this example, a structure variable `dbinfo_t` with members whose names match each column in the `SELECT` result is used to retrieve one result row without putting multiple host variables in the `FETCH` statement.

```
EXEC SQL BEGIN DECLARE SECTION;
typedef struct
{
    int oid;
    char datname[65];
    long long int size;
} dbinfo_t;

dbinfo_t dbval;
EXEC SQL END DECLARE SECTION;

memset(&dbval, 0, sizeof(dbinfo_t));

EXEC SQL DECLARE cur1 CURSOR FOR SELECT oid, datname, pg_database_size(oid) AS size
FROM pg_database;
EXEC SQL OPEN cur1;

/* when end of result set reached, break out of while loop */
EXEC SQL WHENEVER NOT FOUND DO BREAK;
```

```
while (1)
{
    /* Fetch multiple columns into one structure. */
    EXEC SQL FETCH FROM cur1 INTO :dbval;

    /* Print members of the structure. */
    printf("oid=%d, datname=%s, size=%lld\n", dbval.oid, dbval.datname,
dbval.size);
}

EXEC SQL CLOSE cur1;
```

This example shows following result. (The exact values depend on local circumstances.)

```
oid=1, datname=template1, size=4324580
oid=11510, datname=template0, size=4243460
oid=11511, datname=postgres, size=4324580
oid=313780, datname=testdb, size=8183012
```

Structure host variables “absorb” as many columns as the structure as fields. Additional columns can be assigned to other host variables. For example, the above program could also be restructured like this, with the `size` variable outside the structure:

```
EXEC SQL BEGIN DECLARE SECTION;
typedef struct
{
    int oid;
    char datname[65];
} dbinfo_t;

dbinfo_t dbval;
long long int size;
EXEC SQL END DECLARE SECTION;

memset(&dbval, 0, sizeof(dbinfo_t));

EXEC SQL DECLARE cur1 CURSOR FOR SELECT oid, datname, pg_database_size(oid) AS size
FROM pg_database;
EXEC SQL OPEN cur1;

/* when end of result set reached, break out of while loop */
EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    /* Fetch multiple columns into one structure. */
    EXEC SQL FETCH FROM cur1 INTO :dbval, :size;

    /* Print members of the structure. */
    printf("oid=%d, datname=%s, size=%lld\n", dbval.oid, dbval.datname, size);
}

EXEC SQL CLOSE cur1;
```

37.4.4.3.3. Typedefs

Use the `typedef` keyword to map new types to already existing types.

```
EXEC SQL BEGIN DECLARE SECTION;
typedef char mychartype[40];
```

```
typedef long serial_t;
EXEC SQL END DECLARE SECTION;
```

Note that you could also use:

```
EXEC SQL TYPE serial_t IS long;
```

This declaration does not need to be part of a declare section; that is, you can also write typedefs as normal C statements.

Any word you declare as a `typedef` cannot be used as an SQL keyword in `EXEC SQL` commands later in the same program. For example, this won't work:

```
EXEC SQL BEGIN DECLARE SECTION;
    typedef int start;
EXEC SQL END DECLARE SECTION;
...
EXEC SQL START TRANSACTION;
```

ECPG will report a syntax error for `START TRANSACTION`, because it no longer recognizes `START` as an SQL keyword, only as a `typedef`. (If you have such a conflict, and renaming the `typedef` seems impractical, you could write the SQL command using [dynamic SQL](#).)

Note

In PostgreSQL releases before v16, use of SQL keywords as `typedef` names was likely to result in syntax errors associated with use of the `typedef` itself, rather than use of the name as an SQL keyword. The new behavior is less likely to cause problems when an existing ECPG application is recompiled in a new PostgreSQL release with new keywords.

37.4.4.3.4. Pointers

You can declare pointers to the most common types. Note however that you cannot use pointers as target variables of queries without auto-allocation. See [Section 37.7](#) for more information on auto-allocation.

```
EXEC SQL BEGIN DECLARE SECTION;
    int    *intp;
    char **charp;
EXEC SQL END DECLARE SECTION;
```

37.4.5. Handling Nonprimitive SQL Data Types

This section contains information on how to handle nonscalar and user-defined SQL-level data types in ECPG applications. Note that this is distinct from the handling of host variables of nonprimitive types, described in the previous section.

37.4.5.1. Arrays

Multi-dimensional SQL-level arrays are not directly supported in ECPG. One-dimensional SQL-level arrays can be mapped into C array host variables and vice-versa. However, when creating a statement `ecpg` does not know the types of the columns, so that it cannot check if a C array is input into a corresponding SQL-level array. When processing the output of an SQL statement, `ecpg` has the necessary information and thus checks if both are arrays.

If a query accesses *elements* of an array separately, then this avoids the use of arrays in ECPG. Then, a host variable with a type that can be mapped to the element type should be used. For example, if a column type is array of `integer`, a host variable of type `int` can be used. Also if the element type is `varchar` or `text`, a host variable of type `char[]` or `VARCHAR[]` can be used.

Here is an example. Assume the following table:

```
CREATE TABLE t3 (
```

```
    ii integer[]
);

testdb=> SELECT * FROM t3;
      ii
-----
{1,2,3,4,5}
(1 row)
```

The following example program retrieves the 4th element of the array and stores it into a host variable of type int:

```
EXEC SQL BEGIN DECLARE SECTION;
int ii;
EXEC SQL END DECLARE SECTION;

EXEC SQL DECLARE cur1 CURSOR FOR SELECT ii[4] FROM t3;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    EXEC SQL FETCH FROM cur1 INTO :ii ;
    printf("ii=%d\n", ii);
}

EXEC SQL CLOSE cur1;
```

This example shows the following result:

```
ii=4
```

To map multiple array elements to the multiple elements in an array type host variables each element of array column and each element of the host variable array have to be managed separately, for example:

```
EXEC SQL BEGIN DECLARE SECTION;
int ii_a[8];
EXEC SQL END DECLARE SECTION;

EXEC SQL DECLARE cur1 CURSOR FOR SELECT ii[1], ii[2], ii[3], ii[4] FROM t3;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    EXEC SQL FETCH FROM cur1 INTO :ii_a[0], :ii_a[1], :ii_a[2], :ii_a[3];
    ...
}
```

Note again that

```
EXEC SQL BEGIN DECLARE SECTION;
int ii_a[8];
EXEC SQL END DECLARE SECTION;

EXEC SQL DECLARE cur1 CURSOR FOR SELECT ii FROM t3;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;
```

```
while (1)
{
    /* WRONG */
    EXEC SQL FETCH FROM cur1 INTO :ii_a;
    ...
}
```

would not work correctly in this case, because you cannot map an array type column to an array host variable directly.

Another workaround is to store arrays in their external string representation in host variables of type `char[]` or `VARCHAR[]`. For more details about this representation, see [Section 8.15.2](#). Note that this means that the array cannot be accessed naturally as an array in the host program (without further processing that parses the text representation).

37.4.5.2. Composite Types

Composite types are not directly supported in ECPG, but an easy workaround is possible. The available workarounds are similar to the ones described for arrays above: Either access each attribute separately or use the external string representation.

For the following examples, assume the following type and table:

```
CREATE TYPE comp_t AS (intval integer, textval varchar(32));
CREATE TABLE t4 (compval comp_t);
INSERT INTO t4 VALUES ( (256, 'Postgres Pro') );
```

The most obvious solution is to access each attribute separately. The following program retrieves data from the example table by selecting each attribute of the type `comp_t` separately:

```
EXEC SQL BEGIN DECLARE SECTION;
int intval;
varchar textval[33];
EXEC SQL END DECLARE SECTION;

/* Put each element of the composite type column in the SELECT list. */
EXEC SQL DECLARE cur1 CURSOR FOR SELECT (compval).intval, (compval).textval FROM t4;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    /* Fetch each element of the composite type column into host variables. */
    EXEC SQL FETCH FROM cur1 INTO :intval, :textval;

    printf("intval=%d, textval=%s\n", intval, textval.arr);
}

EXEC SQL CLOSE cur1;
```

To enhance this example, the host variables to store values in the `FETCH` command can be gathered into one structure. For more details about the host variable in the structure form, see [Section 37.4.4.3.2](#). To switch to the structure, the example can be modified as below. The two host variables, `intval` and `textval`, become members of the `comp_t` structure, and the structure is specified on the `FETCH` command.

```
EXEC SQL BEGIN DECLARE SECTION;
typedef struct
{
    int intval;
    varchar textval[33];
}
```



```
} comp_t;

comp_t compval;
EXEC SQL END DECLARE SECTION;

/* Put each element of the composite type column in the SELECT list. */
EXEC SQL DECLARE cur1 CURSOR FOR SELECT (compval).intval, (compval).textval FROM t4;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    /* Put all values in the SELECT list into one structure. */
    EXEC SQL FETCH FROM cur1 INTO :compval;

    printf("intval=%d, textval=%s\n", compval.intval, compval.textval.arr);
}

EXEC SQL CLOSE cur1;
```

Although a structure is used in the `FETCH` command, the attribute names in the `SELECT` clause are specified one by one. This can be enhanced by using a `*` to ask for all attributes of the composite type value.

```
...
EXEC SQL DECLARE cur1 CURSOR FOR SELECT (compval).* FROM t4;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    /* Put all values in the SELECT list into one structure. */
    EXEC SQL FETCH FROM cur1 INTO :compval;

    printf("intval=%d, textval=%s\n", compval.intval, compval.textval.arr);
}
...
```

This way, composite types can be mapped into structures almost seamlessly, even though ECPG does not understand the composite type itself.

Finally, it is also possible to store composite type values in their external string representation in host variables of type `char[]` or `VARCHAR[]`. But that way, it is not easily possible to access the fields of the value from the host program.

37.4.5.3. User-Defined Base Types

New user-defined base types are not directly supported by ECPG. You can use the external string representation and host variables of type `char[]` or `VARCHAR[]`, and this solution is indeed appropriate and sufficient for many types.

Here is an example using the data type `complex` from the example in [Section 39.13](#). The external string representation of that type is `(%f, %f)`, which is defined in the functions `complex_in()` and `complex_out()` functions in [Section 39.13](#). The following example inserts the complex type values `(1, 1)` and `(3, 3)` into the columns `a` and `b`, and select them from the table after that.

```
EXEC SQL BEGIN DECLARE SECTION;
    varchar a[64];
    varchar b[64];
EXEC SQL END DECLARE SECTION;
```

```
EXEC SQL INSERT INTO test_complex VALUES ('(1,1)', '(3,3)');

EXEC SQL DECLARE cur1 CURSOR FOR SELECT a, b FROM test_complex;
EXEC SQL OPEN cur1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    EXEC SQL FETCH FROM cur1 INTO :a, :b;
    printf("a=%s, b=%s\n", a.arr, b.arr);
}

EXEC SQL CLOSE cur1;
```

This example shows following result:

```
a=(1,1), b=(3,3)
```

Another workaround is avoiding the direct use of the user-defined types in ECPG and instead create a function or cast that converts between the user-defined type and a primitive type that ECPG can handle. Note, however, that type casts, especially implicit ones, should be introduced into the type system very carefully.

For example,

```
CREATE FUNCTION create_complex(r double, i double) RETURNS complex
LANGUAGE SQL
IMMUTABLE
AS $$ SELECT $1 * complex '(1,0)' + $2 * complex '(0,1)' $$;
```

After this definition, the following

```
EXEC SQL BEGIN DECLARE SECTION;
double a, b, c, d;
EXEC SQL END DECLARE SECTION;
```

```
a = 1;
b = 2;
c = 3;
d = 4;
```

```
EXEC SQL INSERT INTO test_complex VALUES (create_complex(:a, :b),
    create_complex(:c, :d));
```

has the same effect as

```
EXEC SQL INSERT INTO test_complex VALUES ('(1,2)', '(3,4)');
```

37.4.6. Indicators

The examples above do not handle null values. In fact, the retrieval examples will raise an error if they fetch a null value from the database. To be able to pass null values to the database or retrieve null values from the database, you need to append a second host variable specification to each host variable that contains data. This second host variable is called the *indicator* and contains a flag that tells whether the datum is null, in which case the value of the real host variable is ignored. Here is an example that handles the retrieval of null values correctly:

```
EXEC SQL BEGIN DECLARE SECTION;
VARCHAR val;
int val_ind;
EXEC SQL END DECLARE SECTION;
```

```
...
```

```
EXEC SQL SELECT b INTO :val :val_ind FROM test1;
```

The indicator variable `val_ind` will be zero if the value was not null, and it will be negative if the value was null. (See [Section 37.16](#) to enable Oracle-specific behavior.)

The indicator has another function: if the indicator value is positive, it means that the value is not null, but it was truncated when it was stored in the host variable.

If the argument `-r no_indicator` is passed to the preprocessor `ecpg`, it works in “no-indicator” mode. In no-indicator mode, if no indicator variable is specified, null values are signaled (on input and output) for character string types as empty string and for integer types as the lowest possible value for type (for example, `INT_MIN` for `int`).

37.5. Dynamic SQL

In many cases, the particular SQL statements that an application has to execute are known at the time the application is written. In some cases, however, the SQL statements are composed at run time or provided by an external source. In these cases you cannot embed the SQL statements directly into the C source code, but there is a facility that allows you to call arbitrary SQL statements that you provide in a string variable.

37.5.1. Executing Statements without a Result Set

The simplest way to execute an arbitrary SQL statement is to use the command `EXECUTE IMMEDIATE`. For example:

```
EXEC SQL BEGIN DECLARE SECTION;
const char *stmt = "CREATE TABLE test1 (...);";
EXEC SQL END DECLARE SECTION;
```

```
EXEC SQL EXECUTE IMMEDIATE :stmt;
```

`EXECUTE IMMEDIATE` can be used for SQL statements that do not return a result set (e.g., `DDL`, `INSERT`, `UPDATE`, `DELETE`). You cannot execute statements that retrieve data (e.g., `SELECT`) this way. The next section describes how to do that.

37.5.2. Executing a Statement with Input Parameters

A more powerful way to execute arbitrary SQL statements is to prepare them once and execute the prepared statement as often as you like. It is also possible to prepare a generalized version of a statement and then execute specific versions of it by substituting parameters. When preparing the statement, write question marks where you want to substitute parameters later. For example:

```
EXEC SQL BEGIN DECLARE SECTION;
const char *stmt = "INSERT INTO test1 VALUES(?, ?);";
EXEC SQL END DECLARE SECTION;
```

```
EXEC SQL PREPARE mystmt FROM :stmt;
```

```
...
```

```
EXEC SQL EXECUTE mystmt USING 42, 'foobar';
```

When you don't need the prepared statement anymore, you should deallocate it:

```
EXEC SQL DEALLOCATE PREPARE name;
```

37.5.3. Executing a Statement with a Result Set

To execute an SQL statement with a single result row, `EXECUTE` can be used. To save the result, add an `INTO` clause.

```
EXEC SQL BEGIN DECLARE SECTION;
const char *stmt = "SELECT a, b, c FROM test1 WHERE a > ?";
int v1, v2;
VARCHAR v3[50];
EXEC SQL END DECLARE SECTION;

EXEC SQL PREPARE mystmt FROM :stmt;
...
EXEC SQL EXECUTE mystmt INTO :v1, :v2, :v3 USING 37;
```

An EXECUTE command can have an INTO clause, a USING clause, both, or neither.

If a query is expected to return more than one result row, a cursor should be used, as in the following example. (See [Section 37.3.2](#) for more details about the cursor.)

```
EXEC SQL BEGIN DECLARE SECTION;
char dbaname[128];
char datname[128];
char *stmt = "SELECT u.username as dbaname, d.datname "
            " FROM pg_database d, pg_user u "
            " WHERE d.datdba = u.usesysid";
EXEC SQL END DECLARE SECTION;

EXEC SQL CONNECT TO testdb AS con1 USER testuser;
EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

EXEC SQL PREPARE stmt1 FROM :stmt;

EXEC SQL DECLARE cursor1 CURSOR FOR stmt1;
EXEC SQL OPEN cursor1;

EXEC SQL WHENEVER NOT FOUND DO BREAK;

while (1)
{
    EXEC SQL FETCH cursor1 INTO :dbaname, :datname;
    printf("dbaname=%s, datname=%s\n", dbaname, datname);
}

EXEC SQL CLOSE cursor1;

EXEC SQL COMMIT;
EXEC SQL DISCONNECT ALL;
```

37.6. pgtypes Library

The pgtypes library maps Postgres Pro database types to C equivalents that can be used in C programs. It also offers functions to do basic calculations with those types within C, i.e., without the help of the Postgres Pro server. See the following example:

```
EXEC SQL BEGIN DECLARE SECTION;
    date date1;
    timestamp ts1, tsout;
    interval iv1;
    char *out;
EXEC SQL END DECLARE SECTION;

PGTYPESdate_today(&date1);
EXEC SQL SELECT started, duration INTO :ts1, :iv1 FROM datetbl WHERE d=:date1;
```

```
PGTYPEStimestamp_add_interval(&ts1, &iv1, &tsout);
out = PGTYPEStimestamp_to_asc(&tsout);
printf("Started + duration: %s\n", out);
PGTYPESchar_free(out);
```

37.6.1. Character Strings

Some functions such as `PGTYPESnumeric_to_asc` return a pointer to a freshly allocated character string. These results should be freed with `PGTYPESchar_free` instead of `free`. (This is important only on Windows, where memory allocation and release sometimes need to be done by the same library.)

37.6.2. The numeric Type

The numeric type offers to do calculations with arbitrary precision. See [Section 8.1](#) for the equivalent type in the Postgres Pro server. Because of the arbitrary precision this variable needs to be able to expand and shrink dynamically. That's why you can only create numeric variables on the heap, by means of the `PGTYPESnumeric_new` and `PGTYPESnumeric_free` functions. The decimal type, which is similar but limited in precision, can be created on the stack as well as on the heap.

The following functions can be used to work with the numeric type:

`PGTYPESnumeric_new`

Request a pointer to a newly allocated numeric variable.

```
numeric *PGTYPESnumeric_new(void);
```

`PGTYPESnumeric_free`

Free a numeric type, release all of its memory.

```
void PGTYPESnumeric_free(numeric *var);
```

`PGTYPESnumeric_from_asc`

Parse a numeric type from its string notation.

```
numeric *PGTYPESnumeric_from_asc(char *str, char **endptr);
```

Valid formats are for example: `-2`, `.794`, `+3.44`, `592.49E07` or `-32.84e-4`. If the value could be parsed successfully, a valid pointer is returned, else the `NULL` pointer. At the moment ECPG always parses the complete string and so it currently does not support to store the address of the first invalid character in `*endptr`. You can safely set `endptr` to `NULL`.

`PGTYPESnumeric_to_asc`

Returns a pointer to a string allocated by `malloc` that contains the string representation of the numeric type `num`.

```
char *PGTYPESnumeric_to_asc(numeric *num, int dscale);
```

The numeric value will be printed with `dscale` decimal digits, with rounding applied if necessary. The result must be freed with `PGTYPESchar_free()`.

`PGTYPESnumeric_add`

Add two numeric variables into a third one.

```
int PGTYPESnumeric_add(numeric *var1, numeric *var2, numeric *result);
```

The function adds the variables `var1` and `var2` into the result variable `result`. The function returns 0 on success and -1 in case of error.

`PGTYPESnumeric_sub`

Subtract two numeric variables and return the result in a third one.

```
int PGTYPESEnumeric_sub(numeric *var1, numeric *var2, numeric *result);
```

The function subtracts the variable `var2` from the variable `var1`. The result of the operation is stored in the variable `result`. The function returns 0 on success and -1 in case of error.

`PGTYPESEnumeric_mul`

Multiply two numeric variables and return the result in a third one.

```
int PGTYPESEnumeric_mul(numeric *var1, numeric *var2, numeric *result);
```

The function multiplies the variables `var1` and `var2`. The result of the operation is stored in the variable `result`. The function returns 0 on success and -1 in case of error.

`PGTYPESEnumeric_div`

Divide two numeric variables and return the result in a third one.

```
int PGTYPESEnumeric_div(numeric *var1, numeric *var2, numeric *result);
```

The function divides the variables `var1` by `var2`. The result of the operation is stored in the variable `result`. The function returns 0 on success and -1 in case of error.

`PGTYPESEnumeric_cmp`

Compare two numeric variables.

```
int PGTYPESEnumeric_cmp(numeric *var1, numeric *var2)
```

This function compares two numeric variables. In case of error, `INT_MAX` is returned. On success, the function returns one of three possible results:

- 1, if `var1` is bigger than `var2`
- -1, if `var1` is smaller than `var2`
- 0, if `var1` and `var2` are equal

`PGTYPESEnumeric_from_int`

Convert an int variable to a numeric variable.

```
int PGTYPESEnumeric_from_int(signed int int_val, numeric *var);
```

This function accepts a variable of type `signed int` and stores it in the numeric variable `var`. Upon success, 0 is returned and -1 in case of a failure.

`PGTYPESEnumeric_from_long`

Convert a long int variable to a numeric variable.

```
int PGTYPESEnumeric_from_long(signed long int long_val, numeric *var);
```

This function accepts a variable of type `signed long int` and stores it in the numeric variable `var`. Upon success, 0 is returned and -1 in case of a failure.

`PGTYPESEnumeric_copy`

Copy over one numeric variable into another one.

```
int PGTYPESEnumeric_copy(numeric *src, numeric *dst);
```

This function copies over the value of the variable that `src` points to into the variable that `dst` points to. It returns 0 on success and -1 if an error occurs.

`PGTYPESEnumeric_from_double`

Convert a variable of type double to a numeric.

```
int PGTYPESEnumeric_from_double(double d, numeric *dst);
```

This function accepts a variable of type `double` and stores the result in the variable that `dst` points to. It returns 0 on success and -1 if an error occurs.

`PGTYPESnumeric_to_double`

Convert a variable of type numeric to double.

```
int PGTYPESnumeric_to_double(numeric *nv, double *dp)
```

The function converts the numeric value from the variable that `nv` points to into the double variable that `dp` points to. It returns 0 on success and -1 if an error occurs, including overflow. On overflow, the global variable `errno` will be set to `PGTYPES_NUM_OVERFLOW` additionally.

`PGTYPESnumeric_to_int`

Convert a variable of type numeric to int.

```
int PGTYPESnumeric_to_int(numeric *nv, int *ip);
```

The function converts the numeric value from the variable that `nv` points to into the integer variable that `ip` points to. It returns 0 on success and -1 if an error occurs, including overflow. On overflow, the global variable `errno` will be set to `PGTYPES_NUM_OVERFLOW` additionally.

`PGTYPESnumeric_to_long`

Convert a variable of type numeric to long.

```
int PGTYPESnumeric_to_long(numeric *nv, long *lp);
```

The function converts the numeric value from the variable that `nv` points to into the long integer variable that `lp` points to. It returns 0 on success and -1 if an error occurs, including overflow and underflow. On overflow, the global variable `errno` will be set to `PGTYPES_NUM_OVERFLOW` and on underflow `errno` will be set to `PGTYPES_NUM_UNDERFLOW`.

`PGTYPESnumeric_to_decimal`

Convert a variable of type numeric to decimal.

```
int PGTYPESnumeric_to_decimal(numeric *src, decimal *dst);
```

The function converts the numeric value from the variable that `src` points to into the decimal variable that `dst` points to. It returns 0 on success and -1 if an error occurs, including overflow. On overflow, the global variable `errno` will be set to `PGTYPES_NUM_OVERFLOW` additionally.

`PGTYPESnumeric_from_decimal`

Convert a variable of type decimal to numeric.

```
int PGTYPESnumeric_from_decimal(decimal *src, numeric *dst);
```

The function converts the decimal value from the variable that `src` points to into the numeric variable that `dst` points to. It returns 0 on success and -1 if an error occurs. Since the decimal type is implemented as a limited version of the numeric type, overflow cannot occur with this conversion.

37.6.3. The date Type

The date type in C enables your programs to deal with data of the SQL type date. See [Section 8.5](#) for the equivalent type in the Postgres Pro server.

The following functions can be used to work with the date type:

`PGTYPESdate_from_timestamp`

Extract the date part from a timestamp.

```
date PGTYPESdate_from_timestamp(timestamp dt);
```

The function receives a timestamp as its only argument and returns the extracted date part from this timestamp.

`PGTYPESdate_from_asc`

Parse a date from its textual representation.

```
date PGTYPESdate_from_asc(char *str, char **endptr);
```

The function receives a C `char*` string `str` and a pointer to a C `char*` string `endptr`. At the moment ECPG always parses the complete string and so it currently does not support to store the address of the first invalid character in `*endptr`. You can safely set `endptr` to `NULL`.

Note that the function always assumes MDY-formatted dates and there is currently no variable to change that within ECPG.

[Table 37.2](#) shows the allowed input formats.

Table 37.2. Valid Input Formats for `PGTYPESdate_from_asc`

Input	Result
January 8, 1999	January 8, 1999
1999-01-08	January 8, 1999
1/8/1999	January 8, 1999
1/18/1999	January 18, 1999
01/02/03	February 1, 2003
1999-Jan-08	January 8, 1999
Jan-08-1999	January 8, 1999
08-Jan-1999	January 8, 1999
99-Jan-08	January 8, 1999
08-Jan-99	January 8, 1999
08-Jan-06	January 8, 2006
Jan-08-99	January 8, 1999
19990108	ISO 8601; January 8, 1999
990108	ISO 8601; January 8, 1999
1999.008	year and day of year
J2451187	Julian day
January 8, 99 BC	year 99 before the Common Era

`PGTYPESdate_to_asc`

Return the textual representation of a date variable.

```
char *PGTYPESdate_to_asc(date dDate);
```

The function receives the date `dDate` as its only parameter. It will output the date in the form 1999-01-18, i.e., in the YYYY-MM-DD format. The result must be freed with `PGTYPESchar_free()`.

`PGTYPESdate_julmdy`

Extract the values for the day, the month and the year from a variable of type date.

```
void PGTYPESdate_julmdy(date d, int *mdy);
```

The function receives the date `d` and a pointer to an array of 3 integer values `mdy`. The variable name indicates the sequential order: `mdy[0]` will be set to contain the number of the month, `mdy[1]` will be set to the value of the day and `mdy[2]` will contain the year.

PGTYPEStdate_mdyjul

Create a date value from an array of 3 integers that specify the day, the month and the year of the date.

```
void PGTYPEStdate_mdyjul(int *mdy, date *jdate);
```

The function receives the array of the 3 integers (*mdy*) as its first argument and as its second argument a pointer to a variable of type *date* that should hold the result of the operation.

PGTYPEStdate_dayofweek

Return a number representing the day of the week for a date value.

```
int PGTYPEStdate_dayofweek(date d);
```

The function receives the date variable *d* as its only argument and returns an integer that indicates the day of the week for this date.

- 0 - Sunday
- 1 - Monday
- 2 - Tuesday
- 3 - Wednesday
- 4 - Thursday
- 5 - Friday
- 6 - Saturday

PGTYPEStdate_today

Get the current date.

```
void PGTYPEStdate_today(date *d);
```

The function receives a pointer to a date variable (*d*) that it sets to the current date.

PGTYPEStdate_fmt_asc

Convert a variable of type *date* to its textual representation using a format mask.

```
int PGTYPEStdate_fmt_asc(date dDate, char *fmtstring, char *outbuf);
```

The function receives the date to convert (*dDate*), the format mask (*fmtstring*) and the string that will hold the textual representation of the date (*outbuf*).

On success, 0 is returned and a negative value if an error occurred.

The following literals are the field specifiers you can use:

- *dd* - The number of the day of the month.
- *mm* - The number of the month of the year.
- *yy* - The number of the year as a two digit number.
- *yyyy* - The number of the year as a four digit number.
- *ddd* - The name of the day (abbreviated).
- *mmm* - The name of the month (abbreviated).

All other characters are copied 1:1 to the output string.

[Table 37.3](#) indicates a few possible formats. This will give you an idea of how to use this function. All output lines are based on the same date: November 23, 1959.

Table 37.3. Valid Input Formats for `PGTYPESdate_fmt_asc`

Format	Result
mmddyy	112359
ddmmyy	231159
yyymmdd	591123
yy/mm/dd	59/11/23
yy mm dd	59 11 23
yy.mm.dd	59.11.23
.mm.yyyy.dd.	.11.1959.23.
mmm. dd, yyyy	Nov. 23, 1959
mmm dd yyyy	Nov 23 1959
yyyy dd mm	1959 23 11
ddd, mmm. dd, yyyy	Mon, Nov. 23, 1959
(ddd) mmm. dd, yyyy	(Mon) Nov. 23, 1959

`PGTYPESdate_defmt_asc`

Use a format mask to convert a C `char*` string to a value of type `date`.

```
int PGTYPESdate_defmt_asc(date *d, char *fmt, char *str);
```

The function receives a pointer to the date value that should hold the result of the operation (`d`), the format mask to use for parsing the date (`fmt`) and the C `char*` string containing the textual representation of the date (`str`). The textual representation is expected to match the format mask. However you do not need to have a 1:1 mapping of the string to the format mask. The function only analyzes the sequential order and looks for the literals `yy` or `yyyy` that indicate the position of the year, `mm` to indicate the position of the month and `dd` to indicate the position of the day.

[Table 37.4](#) indicates a few possible formats. This will give you an idea of how to use this function.

Table 37.4. Valid Input Formats for `rdefmtdate`

Format	String	Result
ddmmyy	21-2-54	1954-02-21
ddmmyy	2-12-54	1954-12-02
ddmmyy	20111954	1954-11-20
ddmmyy	130464	1964-04-13
mmm.dd.yyyy	MAR-12-1967	1967-03-12
yy/mm/dd	1954, February 3rd	1954-02-03
mmm.dd.yyyy	041269	1969-04-12
yy/mm/dd	In the year 2525, in the month of July, mankind will be alive on the 28th day	2525-07-28
dd-mm-yy	I said on the 28th of July in the year 2525	2525-07-28
mmm.dd.yyyy	9/14/58	1958-09-14
yy/mm/dd	47/03/29	1947-03-29
mmm.dd.yyyy	oct 28 1975	1975-10-28

Format	String	Result
mmddyy	Nov 14th, 1985	1985-11-14

37.6.4. The timestamp Type

The timestamp type in C enables your programs to deal with data of the SQL type timestamp. See [Section 8.5](#) for the equivalent type in the Postgres Pro server.

The following functions can be used to work with the timestamp type:

`PGTYPEStimestamp_from_asc`

Parse a timestamp from its textual representation into a timestamp variable.

```
timestamp PGTYPEStimestamp_from_asc(char *str, char **endptr);
```

The function receives the string to parse (`str`) and a pointer to a C `char*` (`endptr`). At the moment ECPG always parses the complete string and so it currently does not support to store the address of the first invalid character in `*endptr`. You can safely set `endptr` to `NULL`.

The function returns the parsed timestamp on success. On error, `PGTYPEStimestamp` is returned and `errno` is set to `PGTYPES_TS_BAD_TIMESTAMP`. See [PGTYPEStimestamp](#) for important notes on this value.

In general, the input string can contain any combination of an allowed date specification, a white-space character and an allowed time specification. Note that time zones are not supported by ECPG. It can parse them but does not apply any calculation as the Postgres Pro server does for example. Timezone specifiers are silently discarded.

[Table 37.5](#) contains a few examples for input strings.

Table 37.5. Valid Input Formats for `PGTYPEStimestamp_from_asc`

Input	Result
1999-01-08 04:05:06	1999-01-08 04:05:06
January 8 04:05:06 1999 PST	1999-01-08 04:05:06
1999-Jan-08 04:05:06.789-8	1999-01-08 04:05:06.789 (time zone specifier ignored)
J2451187 04:05-08:00	1999-01-08 04:05:00 (time zone specifier ignored)

`PGTYPEStimestamp_to_asc`

Converts a date to a C `char*` string.

```
char *PGTYPEStimestamp_to_asc(timestamp tstamp);
```

The function receives the timestamp `tstamp` as its only argument and returns an allocated string that contains the textual representation of the timestamp. The result must be freed with `PGTYPEStimestamp_free()`.

`PGTYPEStimestamp_current`

Retrieve the current timestamp.

```
void PGTYPEStimestamp_current(timestamp *ts);
```

The function retrieves the current timestamp and saves it into the timestamp variable that `ts` points to.

`PGTYPEStimestamp_fmt_asc`

Convert a timestamp variable to a C `char*` using a format mask.

```
int PGTYPEStimestamp_fmt_asc(timestamp *ts, char *output, int str_len, char
*fmtstr);
```

The function receives a pointer to the timestamp to convert as its first argument (`ts`), a pointer to the output buffer (`output`), the maximal length that has been allocated for the output buffer (`str_len`) and the format mask to use for the conversion (`fmtstr`).

Upon success, the function returns 0 and a negative value if an error occurred.

You can use the following format specifiers for the format mask. The format specifiers are the same ones that are used in the `strftime` function in `libc`. Any non-format specifier will be copied into the output buffer.

- `%A` - is replaced by national representation of the full weekday name.
- `%a` - is replaced by national representation of the abbreviated weekday name.
- `%B` - is replaced by national representation of the full month name.
- `%b` - is replaced by national representation of the abbreviated month name.
- `%C` - is replaced by (year / 100) as decimal number; single digits are preceded by a zero.
- `%c` - is replaced by national representation of time and date.
- `%D` - is equivalent to `%m/%d/%y`.
- `%d` - is replaced by the day of the month as a decimal number (01-31).
- `%E* %O*` - POSIX locale extensions. The sequences `%Ec %EC %Ex %EX %Ey %EY %Od %Oe %OH %OI %Om %OM %OS %Ou %OU %OV %Ow %OW %Oy` are supposed to provide alternative representations.

Additionally `%OB` implemented to represent alternative months names (used standalone, without day mentioned).

- `%e` - is replaced by the day of month as a decimal number (1-31); single digits are preceded by a blank.
- `%F` - is equivalent to `%Y-%m-%d`.
- `%G` - is replaced by a year as a decimal number with century. This year is the one that contains the greater part of the week (Monday as the first day of the week).
- `%g` - is replaced by the same year as in `%G`, but as a decimal number without century (00-99).
- `%H` - is replaced by the hour (24-hour clock) as a decimal number (00-23).
- `%h` - the same as `%b`.
- `%I` - is replaced by the hour (12-hour clock) as a decimal number (01-12).
- `%j` - is replaced by the day of the year as a decimal number (001-366).
- `%k` - is replaced by the hour (24-hour clock) as a decimal number (0-23); single digits are preceded by a blank.
- `%l` - is replaced by the hour (12-hour clock) as a decimal number (1-12); single digits are preceded by a blank.
- `%M` - is replaced by the minute as a decimal number (00-59).
- `%m` - is replaced by the month as a decimal number (01-12).
- `%n` - is replaced by a newline.
- `%O*` - the same as `%E*`.
- `%p` - is replaced by national representation of either “ante meridiem” or “post meridiem” as appropriate.

- `%R` - is equivalent to `%H:%M`.
- `%r` - is equivalent to `%I:%M:%S %p`.
- `%S` - is replaced by the second as a decimal number (00-60).
- `%s` - is replaced by the number of seconds since the Epoch, UTC.
- `%T` - is equivalent to `%H:%M:%S`
- `%t` - is replaced by a tab.
- `%U` - is replaced by the week number of the year (Sunday as the first day of the week) as a decimal number (00-53).
- `%u` - is replaced by the weekday (Monday as the first day of the week) as a decimal number (1-7).
- `%V` - is replaced by the week number of the year (Monday as the first day of the week) as a decimal number (01-53). If the week containing January 1 has four or more days in the new year, then it is week 1; otherwise it is the last week of the previous year, and the next week is week 1.
- `%v` - is equivalent to `%e-%b-%Y`.
- `%W` - is replaced by the week number of the year (Monday as the first day of the week) as a decimal number (00-53).
- `%w` - is replaced by the weekday (Sunday as the first day of the week) as a decimal number (0-6).
- `%X` - is replaced by national representation of the time.
- `%x` - is replaced by national representation of the date.
- `%Y` - is replaced by the year with century as a decimal number.
- `%y` - is replaced by the year without century as a decimal number (00-99).
- `%Z` - is replaced by the time zone name.
- `%z` - is replaced by the time zone offset from UTC; a leading plus sign stands for east of UTC, a minus sign for west of UTC, hours and minutes follow with two digits each and no delimiter between them (common form for [RFC 822](#) date headers).
- `%+` - is replaced by national representation of the date and time.
- `%-*` - GNU libc extension. Do not do any padding when performing numerical outputs.
- `$_*` - GNU libc extension. Explicitly specify space for padding.
- `%0*` - GNU libc extension. Explicitly specify zero for padding.
- `%%` - is replaced by `%`.

`PGTYPEStimestamp_sub`

Subtract one timestamp from another one and save the result in a variable of type interval.

```
int PGTYPEStimestamp_sub(timestamp *ts1, timestamp *ts2, interval *iv);
```

The function will subtract the timestamp variable that `ts2` points to from the timestamp variable that `ts1` points to and will store the result in the interval variable that `iv` points to.

Upon success, the function returns 0 and a negative value if an error occurred.

`PGTYPEStimestamp_defmt_asc`

Parse a timestamp value from its textual representation using a formatting mask.

```
int PGTYPEStimestamp_defmt_asc(char *str, char *fmt, timestamp *d);
```

The function receives the textual representation of a timestamp in the variable `str` as well as the formatting mask to use in the variable `fmt`. The result will be stored in the variable that `d` points to.

If the formatting mask `fmt` is NULL, the function will fall back to the default formatting mask which is `%Y-%m-%d %H:%M:%S`.

This is the reverse function to [PGTYPEStimestamp_fmt_asc](#). See the documentation there in order to find out about the possible formatting mask entries.

`PGTYPEStimestamp_add_interval`

Add an interval variable to a timestamp variable.

```
int PGTYPEStimestamp_add_interval(timestamp *tin, interval *span, timestamp *tout);
```

The function receives a pointer to a timestamp variable `tin` and a pointer to an interval variable `span`. It adds the interval to the timestamp and saves the resulting timestamp in the variable that `tout` points to.

Upon success, the function returns 0 and a negative value if an error occurred.

`PGTYPEStimestamp_sub_interval`

Subtract an interval variable from a timestamp variable.

```
int PGTYPEStimestamp_sub_interval(timestamp *tin, interval *span, timestamp *tout);
```

The function subtracts the interval variable that `span` points to from the timestamp variable that `tin` points to and saves the result into the variable that `tout` points to.

Upon success, the function returns 0 and a negative value if an error occurred.

37.6.5. The interval Type

The interval type in C enables your programs to deal with data of the SQL type interval. See [Section 8.5](#) for the equivalent type in the Postgres Pro server.

The following functions can be used to work with the interval type:

`PGTYPEStimestamp_new`

Return a pointer to a newly allocated interval variable.

```
interval *PGTYPEStimestamp_new(void);
```

`PGTYPEStimestamp_free`

Release the memory of a previously allocated interval variable.

```
void PGTYPEStimestamp_free(interval *intvl);
```

`PGTYPEStimestamp_from_asc`

Parse an interval from its textual representation.

```
interval *PGTYPEStimestamp_from_asc(char *str, char **endptr);
```

The function parses the input string `str` and returns a pointer to an allocated interval variable. At the moment ECPG always parses the complete string and so it currently does not support to store the address of the first invalid character in `*endptr`. You can safely set `endptr` to NULL.

`PGTYPEStimestamp_to_asc`

Convert a variable of type interval to its textual representation.

```
char *PGTYPEStimestamp_to_asc(interval *span);
```

The function converts the interval variable that `span` points to into a C `char*`. The output looks like this example: @ 1 day 12 hours 59 mins 10 secs. The result must be freed with `PGTYPESchar_free()`.

`PGTYPESinterval_copy`

Copy a variable of type interval.

```
int PGTYPESinterval_copy(interval *intvlsrc, interval *intvldest);
```

The function copies the interval variable that `intvlsrc` points to into the variable that `intvldest` points to. Note that you need to allocate the memory for the destination variable before.

37.6.6. The decimal Type

The decimal type is similar to the numeric type. However it is limited to a maximum precision of 30 significant digits. In contrast to the numeric type which can be created on the heap only, the decimal type can be created either on the stack or on the heap (by means of the functions `PGTYPESdecimal_new` and `PGTYPESdecimal_free`). There are a lot of other functions that deal with the decimal type in the Informix compatibility mode described in [Section 37.15](#).

The following functions can be used to work with the decimal type and are not only contained in the `libcompat` library.

`PGTYPESdecimal_new`

Request a pointer to a newly allocated decimal variable.

```
decimal *PGTYPESdecimal_new(void);
```

`PGTYPESdecimal_free`

Free a decimal type, release all of its memory.

```
void PGTYPESdecimal_free(decimal *var);
```

37.6.7. errno Values of pgtypeslib

`PGTYPES_NUM_BAD_NUMERIC`

An argument should contain a numeric variable (or point to a numeric variable) but in fact its in-memory representation was invalid.

`PGTYPES_NUM_OVERFLOW`

An overflow occurred. Since the numeric type can deal with almost arbitrary precision, converting a numeric variable into other types might cause overflow.

`PGTYPES_NUM_UNDERFLOW`

An underflow occurred. Since the numeric type can deal with almost arbitrary precision, converting a numeric variable into other types might cause underflow.

`PGTYPES_NUM_DIVIDE_ZERO`

A division by zero has been attempted.

`PGTYPES_DATE_BAD_DATE`

An invalid date string was passed to the `PGTYPESdate_from_asc` function.

`PGTYPES_DATE_ERR_EARGS`

Invalid arguments were passed to the `PGTYPESdate_defmt_asc` function.

`PGTYPES_DATE_ERR_ENOSHORTDATE`

An invalid token in the input string was found by the `PGTYPESdate_defmt_asc` function.

`PGTYPES_INTERVAL_BAD_INTERVAL`

An invalid interval string was passed to the `PGTYPESinterval_from_asc` function, or an invalid interval value was passed to the `PGTYPESinterval_to_asc` function.

`PGTYPES_DATE_ERR_ENOTDMY`

There was a mismatch in the day/month/year assignment in the `PGTYPESdate_defmt_asc` function.

`PGTYPES_DATE_BAD_DAY`

An invalid day of the month value was found by the `PGTYPESdate_defmt_asc` function.

`PGTYPES_DATE_BAD_MONTH`

An invalid month value was found by the `PGTYPESdate_defmt_asc` function.

`PGTYPES_TS_BAD_TIMESTAMP`

An invalid timestamp string was passed to the `PGTYPEStimestamp_from_asc` function, or an invalid timestamp value was passed to the `PGTYPEStimestamp_to_asc` function.

`PGTYPES_TS_ERR_EINFTIME`

An infinite timestamp value was encountered in a context that cannot handle it.

37.6.8. Special Constants of `pgtypeslib`

`PGTYPESInvalidTimestamp`

A value of type timestamp representing an invalid time stamp. This is returned by the function `PGTYPEStimestamp_from_asc` on parse error. Note that due to the internal representation of the timestamp data type, `PGTYPESInvalidTimestamp` is also a valid timestamp at the same time. It is set to 1899-12-31 23:59:59. In order to detect errors, make sure that your application does not only test for `PGTYPESInvalidTimestamp` but also for `errno != 0` after each call to `PGTYPEStimestamp_from_asc`.

37.7. Using Descriptor Areas

An SQL descriptor area is a more sophisticated method for processing the result of a `SELECT`, `FETCH` or a `DESCRIBE` statement. An SQL descriptor area groups the data of one row of data together with metadata items into one data structure. The metadata is particularly useful when executing dynamic SQL statements, where the nature of the result columns might not be known ahead of time. Postgres Pro provides two ways to use Descriptor Areas: the named SQL Descriptor Areas and the C-structure SQLDAs.

37.7.1. Named SQL Descriptor Areas

A named SQL descriptor area consists of a header, which contains information concerning the entire descriptor, and one or more item descriptor areas, which basically each describe one column in the result row.

Before you can use an SQL descriptor area, you need to allocate one:

```
EXEC SQL ALLOCATE DESCRIPTOR identifier;
```

The identifier serves as the “variable name” of the descriptor area. When you don't need the descriptor anymore, you should deallocate it:

```
EXEC SQL DEALLOCATE DESCRIPTOR identifier;
```

To use a descriptor area, specify it as the storage target in an `INTO` clause, instead of listing host variables:

```
EXEC SQL FETCH NEXT FROM mycursor INTO SQL DESCRIPTOR mydesc;
```


If the result set is empty, the Descriptor Area will still contain the metadata from the query, i.e., the field names.

For not yet executed prepared queries, the `DESCRIBE` statement can be used to get the metadata of the result set:

```
EXEC SQL BEGIN DECLARE SECTION;
char *sql_stmt = "SELECT * FROM table1";
EXEC SQL END DECLARE SECTION;

EXEC SQL PREPARE stmt1 FROM :sql_stmt;
EXEC SQL DESCRIBE stmt1 INTO SQL DESCRIPTOR mydesc;
```

Before PostgreSQL 9.0, the `SQL` keyword was optional, so using `DESCRIPTOR` and `SQL DESCRIPTOR` produced named SQL Descriptor Areas. Now it is mandatory, omitting the `SQL` keyword produces SQLDA Descriptor Areas, see [Section 37.7.2](#).

In `DESCRIBE` and `FETCH` statements, the `INTO` and `USING` keywords can be used to similarly: they produce the result set and the metadata in a Descriptor Area.

Now how do you get the data out of the descriptor area? You can think of the descriptor area as a structure with named fields. To retrieve the value of a field from the header and store it into a host variable, use the following command:

```
EXEC SQL GET DESCRIPTOR name :hostvar = field;
```

Currently, there is only one header field defined: `COUNT`, which tells how many item descriptor areas exist (that is, how many columns are contained in the result). The host variable needs to be of an integer type. To get a field from the item descriptor area, use the following command:

```
EXEC SQL GET DESCRIPTOR name VALUE num :hostvar = field;
```

`num` can be a literal integer or a host variable containing an integer. Possible fields are:

`CARDINALITY` (integer)

number of rows in the result set

`DATA`

actual data item (therefore, the data type of this field depends on the query)

`DATETIME_INTERVAL_CODE` (integer)

When `TYPE` is 9, `DATETIME_INTERVAL_CODE` will have a value of 1 for `DATE`, 2 for `TIME`, 3 for `TIMESTAMP`, 4 for `TIME WITH TIME ZONE`, or 5 for `TIMESTAMP WITH TIME ZONE`.

`DATETIME_INTERVAL_PRECISION` (integer)

not implemented

`INDICATOR` (integer)

the indicator (indicating a null value or a value truncation)

`KEY_MEMBER` (integer)

not implemented

`LENGTH` (integer)

length of the datum in characters

`NAME` (string)

name of the column

NULLABLE (integer)

not implemented

OCTET_LENGTH (integer)

length of the character representation of the datum in bytes

PRECISION (integer)

precision (for type `numeric`)

RETURNED_LENGTH (integer)

length of the datum in characters

RETURNED_OCTET_LENGTH (integer)

length of the character representation of the datum in bytes

SCALE (integer)

scale (for type `numeric`)

TYPE (integer)

numeric code of the data type of the column

In `EXECUTE`, `DECLARE` and `OPEN` statements, the effect of the `INTO` and `USING` keywords are different. A Descriptor Area can also be manually built to provide the input parameters for a query or a cursor and `USING SQL DESCRIPTOR name` is the way to pass the input parameters into a parameterized query. The statement to build a named SQL Descriptor Area is below:

```
EXEC SQL SET DESCRIPTOR name VALUE num field = :hostvar;
```

Postgres Pro supports retrieving more than one record in one `FETCH` statement and storing the data in host variables in this case assumes that the variable is an array. E.g.:

```
EXEC SQL BEGIN DECLARE SECTION;
int id[5];
EXEC SQL END DECLARE SECTION;
```

```
EXEC SQL FETCH 5 FROM mycursor INTO SQL DESCRIPTOR mydesc;
```

```
EXEC SQL GET DESCRIPTOR mydesc VALUE 1 :id = DATA;
```

37.7.2. SQLDA Descriptor Areas

An SQLDA Descriptor Area is a C language structure which can be also used to get the result set and the metadata of a query. One structure stores one record from the result set.

```
EXEC SQL include sqllda.h;
sqllda_t      *mysqlda;
```

```
EXEC SQL FETCH 3 FROM mycursor INTO DESCRIPTOR mysqlda;
```

Note that the `SQL` keyword is omitted. The paragraphs about the use cases of the `INTO` and `USING` keywords in [Section 37.7.1](#) also apply here with an addition. In a `DESCRIBE` statement the `DESCRIPTOR` keyword can be completely omitted if the `INTO` keyword is used:

```
EXEC SQL DESCRIBE prepared_statement INTO mysqlda;
```

The general flow of a program that uses SQLDA is:

1. Prepare a query, and declare a cursor for it.

2. Declare an SQLDA for the result rows.
3. Declare an SQLDA for the input parameters, and initialize them (memory allocation, parameter settings).
4. Open a cursor with the input SQLDA.
5. Fetch rows from the cursor, and store them into an output SQLDA.
6. Read values from the output SQLDA into the host variables (with conversion if necessary).
7. Close the cursor.
8. Free the memory area allocated for the input SQLDA.

37.7.2.1. SQLDA Data Structure

SQLDA uses three data structure types: `sqllda_t`, `sqlvar_t`, and `struct sqlname`.

Tip

Postgres Pro's SQLDA has a similar data structure to the one in IBM DB2 Universal Database, so some technical information on DB2's SQLDA could help understanding Postgres Pro's one better.

37.7.2.1.1. `sqllda_t` Structure

The structure type `sqllda_t` is the type of the actual SQLDA. It holds one record. And two or more `sqllda_t` structures can be connected in a linked list with the pointer in the `desc_next` field, thus representing an ordered collection of rows. So, when two or more rows are fetched, the application can read them by following the `desc_next` pointer in each `sqllda_t` node.

The definition of `sqllda_t` is:

```
struct sqllda_struct
{
    char            sqldaid[8];
    long            sqldabc;
    short           sqln;
    short           sqld;
    struct sqllda_struct *desc_next;
    struct sqlvar_struct sqlvar[1];
};

typedef struct sqllda_struct sqllda_t;
```

The meaning of the fields is:

`sqldaid`

It contains the literal string "SQLDA".

`sqldabc`

It contains the size of the allocated space in bytes.

`sqln`

It contains the number of input parameters for a parameterized query in case it's passed into `OPEN`, `DECLARE` or `EXECUTE` statements using the `USING` keyword. In case it's used as output of `SELECT`, `EXECUTE` or `FETCH` statements, its value is the same as `sqld` statement

`sqld`

It contains the number of fields in a result set.

`desc_next`

If the query returns more than one record, multiple linked SQLDA structures are returned, and `desc_next` holds a pointer to the next entry in the list.

`sqlvar`

This is the array of the columns in the result set.

37.7.2.1.2. `sqlvar_t` Structure

The structure type `sqlvar_t` holds a column value and metadata such as type and length. The definition of the type is:

```
struct sqlvar_struct
{
    short          sqltype;
    short          sqllen;
    char           *sqldata;
    short          *sqlind;
    struct sqlname sqlname;
};

typedef struct sqlvar_struct sqlvar_t;
```

The meaning of the fields is:

`sqltype`

Contains the type identifier of the field. For values, see `enum ECPGttype` in `ecpgtype.h`.

`sqllen`

Contains the binary length of the field. e.g., 4 bytes for `ECPGt_int`.

`sqldata`

Points to the data. The format of the data is described in [Section 37.4.4](#).

`sqlind`

Points to the null indicator. 0 means not null, -1 means null.

`sqlname`

The name of the field.

37.7.2.1.3. `struct sqlname` Structure

A `struct sqlname` structure holds a column name. It is used as a member of the `sqlvar_t` structure. The definition of the structure is:

```
#define NAMEDATALEN 64

struct sqlname
{
    short          length;
    char           data[NAMEDATALEN];
};
```

The meaning of the fields is:

`length`

Contains the length of the field name.

`data`

Contains the actual field name.

37.7.2.2. Retrieving a Result Set Using an SQLDA

The general steps to retrieve a query result set through an SQLDA are:

1. Declare an `sqllda_t` structure to receive the result set.
2. Execute `FETCH/EXECUTE/DESCRIBE` commands to process a query specifying the declared SQLDA.
3. Check the number of records in the result set by looking at `sqln`, a member of the `sqllda_t` structure.
4. Get the values of each column from `sqlvar[0]`, `sqlvar[1]`, etc., members of the `sqllda_t` structure.
5. Go to next row (`sqllda_t` structure) by following the `desc_next` pointer, a member of the `sqllda_t` structure.
6. Repeat above as you need.

Here is an example retrieving a result set through an SQLDA.

First, declare a `sqllda_t` structure to receive the result set.

```
sqllda_t *sqllda1;
```

Next, specify the SQLDA in a command. This is a `FETCH` command example.

```
EXEC SQL FETCH NEXT FROM cur1 INTO DESCRIPTOR sqllda1;
```

Run a loop following the linked list to retrieve the rows.

```
sqllda_t *cur_sqllda;  
  
for (cur_sqllda = sqllda1;  
     cur_sqllda != NULL;  
     cur_sqllda = cur_sqllda->desc_next)  
{  
    ...  
}
```

Inside the loop, run another loop to retrieve each column data (`sqlvar_t` structure) of the row.

```
for (i = 0; i < cur_sqllda->sqln; i++)  
{  
    sqlvar_t v = cur_sqllda->sqlvar[i];  
    char *sqldata = v.sqldata;  
    short sqllen = v.sqllen;  
    ...  
}
```

To get a column value, check the `sqltype` value, a member of the `sqlvar_t` structure. Then, switch to an appropriate way, depending on the column type, to copy data from the `sqlvar` field to a host variable.

```
char var_buf[1024];  
  
switch (v.sqltype)  
{  
    case ECPGt_char:  
        memset(&var_buf, 0, sizeof(var_buf));  
        memcpy(&var_buf, sqldata, (sizeof(var_buf) <= sqllen ? sizeof(var_buf) - 1 :  
sqllen));  
        break;  
  
    case ECPGt_int: /* integer */  
        memcpy(&intval, sqldata, sqllen);  
        snprintf(var_buf, sizeof(var_buf), "%d", intval);  
}
```

```
        break;

        ...
    }
```

37.7.2.3. Passing Query Parameters Using an SQLDA

The general steps to use an SQLDA to pass input parameters to a prepared query are:

1. Create a prepared query (prepared statement)
2. Declare an `sqlda_t` structure as an input SQLDA.
3. Allocate memory area (as `sqlda_t` structure) for the input SQLDA.
4. Set (copy) input values in the allocated memory.
5. Open a cursor with specifying the input SQLDA.

Here is an example.

First, create a prepared statement.

```
EXEC SQL BEGIN DECLARE SECTION;
char query[1024] = "SELECT d.oid, * FROM pg_database d, pg_stat_database s WHERE d.oid
= s.datid AND (d.datname = ? OR d.oid = ?)";
EXEC SQL END DECLARE SECTION;

EXEC SQL PREPARE stmt1 FROM :query;
```

Next, allocate memory for an SQLDA, and set the number of input parameters in `sqln`, a member variable of the `sqlda_t` structure. When two or more input parameters are required for the prepared query, the application has to allocate additional memory space which is calculated by $(\text{nr. of params} - 1) * \text{sizeof}(\text{sqlvar_t})$. The example shown here allocates memory space for two input parameters.

```
sqlda_t *sqlda2;

sqlda2 = (sqlda_t *) malloc(sizeof(sqlda_t) + sizeof(sqlvar_t));
memset(sqlda2, 0, sizeof(sqlda_t) + sizeof(sqlvar_t));

sqlda2->sqln = 2; /* number of input variables */
```

After memory allocation, store the parameter values into the `sqlvar[]` array. (This is same array used for retrieving column values when the SQLDA is receiving a result set.) In this example, the input parameters are "postgres", having a string type, and 1, having an integer type.

```
sqlda2->sqlvar[0].sqltype = ECPGt_char;
sqlda2->sqlvar[0].sqldata = "postgres";
sqlda2->sqlvar[0].sqlllen = 8;

int intval = 1;
sqlda2->sqlvar[1].sqltype = ECPGt_int;
sqlda2->sqlvar[1].sqldata = (char *) &intval;
sqlda2->sqlvar[1].sqlllen = sizeof(intval);
```

By opening a cursor and specifying the SQLDA that was set up beforehand, the input parameters are passed to the prepared statement.

```
EXEC SQL OPEN cur1 USING DESCRIPTOR sqlda2;
```

Finally, after using input SQLDAs, the allocated memory space must be freed explicitly, unlike SQLDAs used for receiving query results.

```
free(sqlda2);
```

37.7.2.4. A Sample Application Using SQLDA

Here is an example program, which describes how to fetch access statistics of the databases, specified by the input parameters, from the system catalogs.

This application joins two system tables, `pg_database` and `pg_stat_database` on the database OID, and also fetches and shows the database statistics which are retrieved by two input parameters (a database `postgres`, and OID 1).

First, declare an SQLDA for input and an SQLDA for output.

```
EXEC SQL include sqllda.h;
```

```
sqllda_t *sqllda1; /* an output descriptor */
sqllda_t *sqllda2; /* an input descriptor  */
```

Next, connect to the database, prepare a statement, and declare a cursor for the prepared statement.

```
int
main(void)
{
    EXEC SQL BEGIN DECLARE SECTION;
    char query[1024] = "SELECT d.oid,* FROM pg_database d, pg_stat_database s WHERE
d.oid=s.datid AND ( d.datname=? OR d.oid=? )";
    EXEC SQL END DECLARE SECTION;

    EXEC SQL CONNECT TO testdb AS con1 USER testuser;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    EXEC SQL PREPARE stmt1 FROM :query;
    EXEC SQL DECLARE cur1 CURSOR FOR stmt1;
```

Next, put some values in the input SQLDA for the input parameters. Allocate memory for the input SQLDA, and set the number of input parameters to `sqln`. Store type, value, and value length into `sqltype`, `sqldata`, and `sqlllen` in the `sqlvar` structure.

```
/* Create SQLDA structure for input parameters. */
sqllda2 = (sqllda_t *) malloc(sizeof(sqllda_t) + sizeof(sqlvar_t));
memset(sqllda2, 0, sizeof(sqllda_t) + sizeof(sqlvar_t));
sqllda2->sqln = 2; /* number of input variables */

sqllda2->sqlvar[0].sqltype = ECPGt_char;
sqllda2->sqlvar[0].sqldata = "postgres";
sqllda2->sqlvar[0].sqlllen = 8;

intval = 1;
sqllda2->sqlvar[1].sqltype = ECPGt_int;
sqllda2->sqlvar[1].sqldata = (char *)&intval;
sqllda2->sqlvar[1].sqlllen = sizeof(intval);
```

After setting up the input SQLDA, open a cursor with the input SQLDA.

```
/* Open a cursor with input parameters. */
EXEC SQL OPEN cur1 USING DESCRIPTOR sqllda2;
```

Fetch rows into the output SQLDA from the opened cursor. (Generally, you have to call `FETCH` repeatedly in the loop, to fetch all rows in the result set.)

```
while (1)
{
    sqllda_t *cur_sqllda;
```

```

/* Assign descriptor to the cursor */
EXEC SQL FETCH NEXT FROM cur1 INTO DESCRIPTOR sqlda1;

```

Next, retrieve the fetched records from the SQLDA, by following the linked list of the `sqlda_t` structure.

```

for (cur_sqlda = sqlda1 ;
     cur_sqlda != NULL ;
     cur_sqlda = cur_sqlda->desc_next)
{
    ...
}

```

Read each columns in the first record. The number of columns is stored in `sqld`, the actual data of the first column is stored in `sqlvar[0]`, both members of the `sqlda_t` structure.

```

/* Print every column in a row. */
for (i = 0; i < sqlda1->sqld; i++)
{
    sqlvar_t v = sqlda1->sqlvar[i];
    char *sqldata = v.sqldata;
    short sqllen = v.sqlllen;

    strncpy(name_buf, v.sqlname.data, v.sqlname.length);
    name_buf[v.sqlname.length] = '\0';
}

```

Now, the column data is stored in the variable `v`. Copy every datum into host variables, looking at `v.sqltype` for the type of the column.

```

switch (v.sqltype) {
    int intval;
    double doubleval;
    unsigned long long int longlongval;

    case ECPGt_char:
        memset(&var_buf, 0, sizeof(var_buf));
        memcpy(&var_buf, sqldata, (sizeof(var_buf) <= sqllen ?
sizeof(var_buf)-1 : sqllen));
        break;

    case ECPGt_int: /* integer */
        memcpy(&intval, sqldata, sqllen);
        snprintf(var_buf, sizeof(var_buf), "%d", intval);
        break;

    ...

    default:
        ...
}

printf("%s = %s (type: %d)\n", name_buf, var_buf, v.sqltype);
}

```

Close the cursor after processing all of records, and disconnect from the database.

```

EXEC SQL CLOSE cur1;
EXEC SQL COMMIT;

EXEC SQL DISCONNECT ALL;

```

The whole program is shown in [Example 37.1](#).

Example 37.1. Example SQLDA Program

```
#include <stdlib.h>
#include <string.h>
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>

EXEC SQL include sqllda.h;

sqllda_t *sqllda1; /* descriptor for output */
sqllda_t *sqllda2; /* descriptor for input */

EXEC SQL WHENEVER NOT FOUND DO BREAK;
EXEC SQL WHENEVER SQLERROR STOP;

int
main(void)
{
    EXEC SQL BEGIN DECLARE SECTION;
    char query[1024] = "SELECT d.oid,* FROM pg_database d, pg_stat_database s WHERE
d.oid=s.datid AND ( d.datname=? OR d.oid=? )";

    int intval;
    unsigned long long int longlongval;
    EXEC SQL END DECLARE SECTION;

    EXEC SQL CONNECT TO uptimedb AS con1 USER uptime;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    EXEC SQL PREPARE stmt1 FROM :query;
    EXEC SQL DECLARE cur1 CURSOR FOR stmt1;

    /* Create an SQLDA structure for an input parameter */
    sqllda2 = (sqllda_t *)malloc(sizeof(sqllda_t) + sizeof(sqlvar_t));
    memset(sqllda2, 0, sizeof(sqllda_t) + sizeof(sqlvar_t));
    sqllda2->sqln = 2; /* a number of input variables */

    sqllda2->sqlvar[0].sqltype = ECPGt_char;
    sqllda2->sqlvar[0].sqldata = "postgres";
    sqllda2->sqlvar[0].sqlllen = 8;

    intval = 1;
    sqllda2->sqlvar[1].sqltype = ECPGt_int;
    sqllda2->sqlvar[1].sqldata = (char *) &intval;
    sqllda2->sqlvar[1].sqlllen = sizeof(intval);

    /* Open a cursor with input parameters. */
    EXEC SQL OPEN cur1 USING DESCRIPTOR sqllda2;

    while (1)
    {
        sqllda_t *cur_sqllda;

        /* Assign descriptor to the cursor */
        EXEC SQL FETCH NEXT FROM cur1 INTO DESCRIPTOR sqllda1;

        for (cur_sqllda = sqllda1 ;
```

```
    cur_sqlda != NULL ;
    cur_sqlda = cur_sqlda->desc_next)
{
    int i;
    char name_buf[1024];
    char var_buf[1024];

    /* Print every column in a row. */
    for (i=0 ; i<cur_sqlda->sqld ; i++)
    {
        sqlvar_t v = cur_sqlda->sqlvar[i];
        char *sqldata = v.sqldata;
        short sqllen = v.sqlllen;

        strncpy(name_buf, v.sqlname.data, v.sqlname.length);
        name_buf[v.sqlname.length] = '\0';

        switch (v.sqltype)
        {
            case ECPGt_char:
                memset(&var_buf, 0, sizeof(var_buf));
                memcpy(&var_buf, sqldata, (sizeof(var_buf)<=sqllen ?
sizeof(var_buf)-1 : sqllen) );
                break;

            case ECPGt_int: /* integer */
                memcpy(&intval, sqldata, sqllen);
                snprintf(var_buf, sizeof(var_buf), "%d", intval);
                break;

            case ECPGt_long_long: /* bigint */
                memcpy(&longlongval, sqldata, sqllen);
                snprintf(var_buf, sizeof(var_buf), "%lld", longlongval);
                break;

            default:
            {
                int i;
                memset(var_buf, 0, sizeof(var_buf));
                for (i = 0; i < sqllen; i++)
                {
                    char tmpbuf[16];
                    snprintf(tmpbuf, sizeof(tmpbuf), "%02x ", (unsigned char)
sqldata[i]);
                    strncat(var_buf, tmpbuf, sizeof(var_buf));
                }
                break;
            }

            printf("%s = %s (type: %d)\n", name_buf, var_buf, v.sqltype);
        }

        printf("\n");
    }
}

EXEC SQL CLOSE cur1;
```

```
EXEC SQL COMMIT;

EXEC SQL DISCONNECT ALL;

return 0;
}
```

The output of this example should look something like the following (some numbers will vary).

```
oid = 1 (type: 1)
datname = template1 (type: 1)
datdba = 10 (type: 1)
encoding = 0 (type: 5)
datistemplate = t (type: 1)
dataallowconn = t (type: 1)
dathasloginevt = f (type: 1)
datconnlimit = -1 (type: 5)
datfrozenxid = 379 (type: 1)
dattablespace = 1663 (type: 1)
datconfig = (type: 1)
datacl = {=c/uptime,uptime=CTc/uptime} (type: 1)
datid = 1 (type: 1)
datname = template1 (type: 1)
numbackends = 0 (type: 5)
xact_commit = 113606 (type: 9)
xact_rollback = 0 (type: 9)
blks_read = 130 (type: 9)
blks_hit = 7341714 (type: 9)
tup_returned = 38262679 (type: 9)
tup_fetched = 1836281 (type: 9)
tup_inserted = 0 (type: 9)
tup_updated = 0 (type: 9)
tup_deleted = 0 (type: 9)
```

```
oid = 11511 (type: 1)
datname = postgres (type: 1)
datdba = 10 (type: 1)
encoding = 0 (type: 5)
datistemplate = f (type: 1)
dataallowconn = t (type: 1)
dathasloginevt = f (type: 1)
datconnlimit = -1 (type: 5)
datfrozenxid = 379 (type: 1)
dattablespace = 1663 (type: 1)
datconfig = (type: 1)
datacl = (type: 1)
datid = 11511 (type: 1)
datname = postgres (type: 1)
numbackends = 0 (type: 5)
xact_commit = 221069 (type: 9)
xact_rollback = 18 (type: 9)
blks_read = 1176 (type: 9)
blks_hit = 13943750 (type: 9)
tup_returned = 77410091 (type: 9)
tup_fetched = 3253694 (type: 9)
tup_inserted = 0 (type: 9)
tup_updated = 0 (type: 9)
tup_deleted = 0 (type: 9)
```

37.8. Error Handling

This section describes how you can handle exceptional conditions and warnings in an embedded SQL program. There are two nonexclusive facilities for this.

- Callbacks can be configured to handle warning and error conditions using the `WHENEVER` command.
- Detailed information about the error or warning can be obtained from the `sqlca` variable.

37.8.1. Setting Callbacks

One simple method to catch errors and warnings is to set a specific action to be executed whenever a particular condition occurs. In general:

```
EXEC SQL WHENEVER condition action;
```

condition can be one of the following:

`SQLERROR`

The specified action is called whenever an error occurs during the execution of an SQL statement.

`SQLWARNING`

The specified action is called whenever a warning occurs during the execution of an SQL statement.

`NOT FOUND`

The specified action is called whenever an SQL statement retrieves or affects zero rows. (This condition is not an error, but you might be interested in handling it specially.)

action can be one of the following:

`CONTINUE`

This effectively means that the condition is ignored. This is the default.

`GOTO label`

`GO TO label`

Jump to the specified label (using a C `goto` statement).

`SQLPRINT`

Print a message to standard error. This is useful for simple programs or during prototyping. The details of the message cannot be configured.

`STOP`

Call `exit(1)`, which will terminate the program.

`DO BREAK`

Execute the C statement `break`. This should only be used in loops or `switch` statements.

`DO CONTINUE`

Execute the C statement `continue`. This should only be used in loops statements. if executed, will cause the flow of control to return to the top of the loop.

`CALL name (args)`

`DO name (args)`

Call the specified C functions with the specified arguments. (This use is different from the meaning of `CALL` and `DO` in the normal Postgres Pro grammar.)

The SQL standard only provides for the actions `CONTINUE` and `GOTO` (and `GO TO`).

Here is an example that you might want to use in a simple program. It prints a simple message when a warning occurs and aborts the program when an error happens:

```
EXEC SQL WHENEVER SQLWARNING SQLPRINT;
EXEC SQL WHENEVER SQLERROR STOP;
```

The statement `EXEC SQL WHENEVER` is a directive of the SQL preprocessor, not a C statement. The error or warning actions that it sets apply to all embedded SQL statements that appear below the point where the handler is set, unless a different action was set for the same condition between the first `EXEC SQL WHENEVER` and the SQL statement causing the condition, regardless of the flow of control in the C program. So neither of the two following C program excerpts will have the desired effect:

```
/*
 * WRONG
 */
int main(int argc, char *argv[])
{
    ...
    if (verbose) {
        EXEC SQL WHENEVER SQLWARNING SQLPRINT;
    }
    ...
    EXEC SQL SELECT ...;
    ...
}

/*
 * WRONG
 */
int main(int argc, char *argv[])
{
    ...
    set_error_handler();
    ...
    EXEC SQL SELECT ...;
    ...
}

static void set_error_handler(void)
{
    EXEC SQL WHENEVER SQLERROR STOP;
}
```

37.8.2. sqlca

For more powerful error handling, the embedded SQL interface provides a global variable with the name `sqlca` (SQL communication area) that has the following structure:

```
struct
{
    char sqlcaid[8];
    long sqlabc;
    long sqlcode;
    struct
    {
        int sqlerrml;
        char sqlerrmc[SQLERRMC_LEN];
    } sqlerrm;
    char sqlerrp[8];
}
```

```
long sqlerrd[6];
char sqlwarn[8];
char sqlstate[5];
} sqlca;
```

(In a multithreaded program, every thread automatically gets its own copy of `sqlca`. This works similarly to the handling of the standard C global variable `errno`.)

`sqlca` covers both warnings and errors. If multiple warnings or errors occur during the execution of a statement, then `sqlca` will only contain information about the last one.

If no error occurred in the last SQL statement, `sqlca.sqlcode` will be 0 and `sqlca.sqlstate` will be "00000". If a warning or error occurred, then `sqlca.sqlcode` will be negative and `sqlca.sqlstate` will be different from "00000". A positive `sqlca.sqlcode` indicates a harmless condition, such as that the last query returned zero rows. `sqlcode` and `sqlstate` are two different error code schemes; details appear below.

If the last SQL statement was successful, then `sqlca.sqlerrd[1]` contains the OID of the processed row, if applicable, and `sqlca.sqlerrd[2]` contains the number of processed or returned rows, if applicable to the command.

In case of an error or warning, `sqlca.sqlerrm.sqlerrmc` will contain a string that describes the error. The field `sqlca.sqlerrm.sqlerrml` contains the length of the error message that is stored in `sqlca.sqlerrm.sqlerrmc` (the result of `strlen()`, not really interesting for a C programmer). Note that some messages are too long to fit in the fixed-size `sqlerrmc` array; they will be truncated.

In case of a warning, `sqlca.sqlwarn[2]` is set to W. (In all other cases, it is set to something different from W.) If `sqlca.sqlwarn[1]` is set to W, then a value was truncated when it was stored in a host variable. `sqlca.sqlwarn[0]` is set to W if any of the other elements are set to indicate a warning.

The fields `sqlcaid`, `sqlabc`, `sqlerrp`, and the remaining elements of `sqlerrd` and `sqlwarn` currently contain no useful information.

The structure `sqlca` is not defined in the SQL standard, but is implemented in several other SQL database systems. The definitions are similar at the core, but if you want to write portable applications, then you should investigate the different implementations carefully.

Here is one example that combines the use of `WHENEVER` and `sqlca`, printing out the contents of `sqlca` when an error occurs. This is perhaps useful for debugging or prototyping applications, before installing a more "user-friendly" error handler.

```
EXEC SQL WHENEVER SQLERROR CALL print_sqlca();
```

```
void
print_sqlca()
{
    fprintf(stderr, "==== sqlca ====\n");
    fprintf(stderr, "sqlcode: %ld\n", sqlca.sqlcode);
    fprintf(stderr, "sqlerrm.sqlerrml: %d\n", sqlca.sqlerrm.sqlerrml);
    fprintf(stderr, "sqlerrm.sqlerrmc: %s\n", sqlca.sqlerrm.sqlerrmc);
    fprintf(stderr, "sqlerrd: %ld %ld %ld %ld %ld %ld\n",
        sqlca.sqlerrd[0], sqlca.sqlerrd[1], sqlca.sqlerrd[2],
        sqlca.sqlerrd[3], sqlca.sqlerrd[4], sqlca.sqlerrd[5]);
    fprintf(stderr, "sqlwarn: %d %d %d %d %d %d %d %d\n", sqlca.sqlwarn[0],
        sqlca.sqlwarn[1], sqlca.sqlwarn[2],
        sqlca.sqlwarn[3],
        sqlca.sqlwarn[4], sqlca.sqlwarn[5],
        sqlca.sqlwarn[6],
        sqlca.sqlwarn[7]);
}
```

```
fprintf(stderr, "sqlstate: %5s\n", sqlca.sqlstate);
fprintf(stderr, "=====\n");
}
```

The result could look as follows (here an error due to a misspelled table name):

```
==== sqlca ====
sqlcode: -400
sqlerrm.sqlerrml: 49
sqlerrm.sqlerrmc: relation "pg_databasep" does not exist on line 38
sqlerrd: 0 0 0 0 0 0
sqlwarn: 0 0 0 0 0 0 0 0
sqlstate: 42P01
=====
```

37.8.3. SQLSTATE vs. SQLCODE

The fields `sqlca.sqlstate` and `sqlca.sqlcode` are two different schemes that provide error codes. Both are derived from the SQL standard, but `SQLCODE` has been marked deprecated in the SQL-92 edition of the standard and has been dropped in later editions. Therefore, new applications are strongly encouraged to use `SQLSTATE`.

`SQLSTATE` is a five-character array. The five characters contain digits or upper-case letters that represent codes of various error and warning conditions. `SQLSTATE` has a hierarchical scheme: the first two characters indicate the general class of the condition, the last three characters indicate a subclass of the general condition. A successful state is indicated by the code `00000`. The `SQLSTATE` codes are for the most part defined in the SQL standard. The Postgres Pro server natively supports `SQLSTATE` error codes; therefore a high degree of consistency can be achieved by using this error code scheme throughout all applications. For further information see [Appendix A](#).

`SQLCODE`, the deprecated error code scheme, is a simple integer. A value of 0 indicates success, a positive value indicates success with additional information, a negative value indicates an error. The SQL standard only defines the positive value +100, which indicates that the last command returned or affected zero rows, and no specific negative values. Therefore, this scheme can only achieve poor portability and does not have a hierarchical code assignment. Historically, the embedded SQL processor for Postgres Pro has assigned some specific `SQLCODE` values for its use, which are listed below with their numeric value and their symbolic name. Remember that these are not portable to other SQL implementations. To simplify the porting of applications to the `SQLSTATE` scheme, the corresponding `SQLSTATE` is also listed. There is, however, no one-to-one or one-to-many mapping between the two schemes (indeed it is many-to-many), so you should consult the global `SQLSTATE` listing in [Appendix A](#) in each case.

These are the assigned `SQLCODE` values:

0 (ECPG_NO_ERROR)

Indicates no error. (SQLSTATE 00000)

100 (ECPG_NOT_FOUND)

This is a harmless condition indicating that the last command retrieved or processed zero rows, or that you are at the end of the cursor. (SQLSTATE 02000)

When processing a cursor in a loop, you could use this code as a way to detect when to abort the loop, like this:

```
while (1)
{
    EXEC SQL FETCH ... ;
    if (sqlca.sqlcode == ECPG_NOT_FOUND)
        break;
}
```

But `WHENEVER NOT FOUND DO BREAK` effectively does this internally, so there is usually no advantage in writing this out explicitly.

-12 (ECPG_OUT_OF_MEMORY)

Indicates that your virtual memory is exhausted. The numeric value is defined as `-ENOMEM`. (SQLSTATE YE001)

-200 (ECPG_UNSUPPORTED)

Indicates the preprocessor has generated something that the library does not know about. Perhaps you are running incompatible versions of the preprocessor and the library. (SQLSTATE YE002)

-201 (ECPG_TOO_MANY_ARGUMENTS)

This means that the command specified more host variables than the command expected. (SQLSTATE 07001 or 07002)

-202 (ECPG_TOO_FEW_ARGUMENTS)

This means that the command specified fewer host variables than the command expected. (SQLSTATE 07001 or 07002)

-203 (ECPG_TOO_MANY_MATCHES)

This means a query has returned multiple rows but the statement was only prepared to store one result row (for example, because the specified variables are not arrays). (SQLSTATE 21000)

-204 (ECPG_INT_FORMAT)

The host variable is of type `int` and the datum in the database is of a different type and contains a value that cannot be interpreted as an `int`. The library uses `strtol()` for this conversion. (SQLSTATE 42804)

-205 (ECPG_UINT_FORMAT)

The host variable is of type `unsigned int` and the datum in the database is of a different type and contains a value that cannot be interpreted as an `unsigned int`. The library uses `strtoul()` for this conversion. (SQLSTATE 42804)

-206 (ECPG_FLOAT_FORMAT)

The host variable is of type `float` and the datum in the database is of another type and contains a value that cannot be interpreted as a `float`. The library uses `strtod()` for this conversion. (SQLSTATE 42804)

-207 (ECPG_NUMERIC_FORMAT)

The host variable is of type `numeric` and the datum in the database is of another type and contains a value that cannot be interpreted as a `numeric` value. (SQLSTATE 42804)

-208 (ECPG_INTERVAL_FORMAT)

The host variable is of type `interval` and the datum in the database is of another type and contains a value that cannot be interpreted as an `interval` value. (SQLSTATE 42804)

-209 (ECPG_DATE_FORMAT)

The host variable is of type `date` and the datum in the database is of another type and contains a value that cannot be interpreted as a `date` value. (SQLSTATE 42804)

-210 (ECPG_TIMESTAMP_FORMAT)

The host variable is of type `timestamp` and the datum in the database is of another type and contains a value that cannot be interpreted as a `timestamp` value. (SQLSTATE 42804)

-211 (ECPG_CONVERT_BOOL)

This means the host variable is of type `bool` and the datum in the database is neither 't' nor 'f'. (SQLSTATE 42804)

-212 (ECPG_EMPTY)

The statement sent to the Postgres Pro server was empty. (This cannot normally happen in an embedded SQL program, so it might point to an internal error.) (SQLSTATE YE002)

-213 (ECPG_MISSING_INDICATOR)

A null value was returned and no null indicator variable was supplied. (SQLSTATE 22002)

-214 (ECPG_NO_ARRAY)

An ordinary variable was used in a place that requires an array. (SQLSTATE 42804)

-215 (ECPG_DATA_NOT_ARRAY)

The database returned an ordinary variable in a place that requires array value. (SQLSTATE 42804)

-216 (ECPG_ARRAY_INSERT)

The value could not be inserted into the array. (SQLSTATE 42804)

-220 (ECPG_NO_CONN)

The program tried to access a connection that does not exist. (SQLSTATE 08003)

-221 (ECPG_NOT_CONN)

The program tried to access a connection that does exist but is not open. (This is an internal error.) (SQLSTATE YE002)

-230 (ECPG_INVALID_STMT)

The statement you are trying to use has not been prepared. (SQLSTATE 26000)

-239 (ECPG_INFORMIX_DUPLICATE_KEY)

Duplicate key error, violation of unique constraint (Informix compatibility mode). (SQLSTATE 23505)

-240 (ECPG_UNKNOWN_DESCRIPTOR)

The descriptor specified was not found. The statement you are trying to use has not been prepared. (SQLSTATE 33000)

-241 (ECPG_INVALID_DESCRIPTOR_INDEX)

The descriptor index specified was out of range. (SQLSTATE 07009)

-242 (ECPG_UNKNOWN_DESCRIPTOR_ITEM)

An invalid descriptor item was requested. (This is an internal error.) (SQLSTATE YE002)

-243 (ECPG_VAR_NOT_NUMERIC)

During the execution of a dynamic statement, the database returned a numeric value and the host variable was not numeric. (SQLSTATE 07006)

-244 (ECPG_VAR_NOT_CHAR)

During the execution of a dynamic statement, the database returned a non-numeric value and the host variable was numeric. (SQLSTATE 07006)

-284 (ECPG_INFORMIX_SUBSELECT_NOT_ONE)

A result of the subquery is not single row (Informix compatibility mode). (SQLSTATE 21000)

-400 (ECPG_PGSQL)

Some error caused by the Postgres Pro server. The message contains the error message from the Postgres Pro server.

-401 (ECPG_TRANS)

The Postgres Pro server signaled that we cannot start, commit, or rollback the transaction. (SQLSTATE 08007)

-402 (ECPG_CONNECT)

The connection attempt to the database did not succeed. (SQLSTATE 08001)

-403 (ECPG_DUPLICATE_KEY)

Duplicate key error, violation of unique constraint. (SQLSTATE 23505)

-404 (ECPG_SUBSELECT_NOT_ONE)

A result for the subquery is not single row. (SQLSTATE 21000)

-602 (ECPG_WARNING_UNKNOWN_PORTAL)

An invalid cursor name was specified. (SQLSTATE 34000)

-603 (ECPG_WARNING_IN_TRANSACTION)

Transaction is in progress. (SQLSTATE 25001)

-604 (ECPG_WARNING_NO_TRANSACTION)

There is no active (in-progress) transaction. (SQLSTATE 25P01)

-605 (ECPG_WARNING_PORTAL_EXISTS)

An existing cursor name was specified. (SQLSTATE 42P03)

37.9. Preprocessor Directives

Several preprocessor directives are available that modify how the `ecpg` preprocessor parses and processes a file.

37.9.1. Including Files

To include an external file into your embedded SQL program, use:

```
EXEC SQL INCLUDE filename;  
EXEC SQL INCLUDE <filename>;  
EXEC SQL INCLUDE "filename";
```

The embedded SQL preprocessor will look for a file named *filename.h*, preprocess it, and include it in the resulting C output. Thus, embedded SQL statements in the included file are handled correctly.

The `ecpg` preprocessor will search a file at several directories in following order:

- current directory
- `/usr/local/include`
- Postgres Pro include directory, defined at build time (e.g., `/usr/local/pgsql/include`)
- `/usr/include`

But when `EXEC SQL INCLUDE "filename"` is used, only the current directory is searched.

In each directory, the preprocessor will first look for the file name as given, and if not found will append `.h` to the file name and try again (unless the specified file name already has that suffix).

Note that `EXEC SQL INCLUDE` is *not* the same as:

```
#include <filename.h>
```

because this file would not be subject to SQL command preprocessing. Naturally, you can continue to use the C `#include` directive to include other header files.

Note

The include file name is case-sensitive, even though the rest of the `EXEC SQL INCLUDE` command follows the normal SQL case-sensitivity rules.

37.9.2. The `define` and `undef` Directives

Similar to the directive `#define` that is known from C, embedded SQL has a similar concept:

```
EXEC SQL DEFINE name;
EXEC SQL DEFINE name value;
```

So you can define a name:

```
EXEC SQL DEFINE HAVE_FEATURE;
```

And you can also define constants:

```
EXEC SQL DEFINE MYNUMBER 12;
EXEC SQL DEFINE MYSTRING 'abc';
```

Use `undef` to remove a previous definition:

```
EXEC SQL UNDEF MYNUMBER;
```

Of course you can continue to use the C versions `#define` and `#undef` in your embedded SQL program. The difference is where your defined values get evaluated. If you use `EXEC SQL DEFINE` then the `ecpg` preprocessor evaluates the defines and substitutes the values. For example if you write:

```
EXEC SQL DEFINE MYNUMBER 12;
...
EXEC SQL UPDATE Tbl SET col = MYNUMBER;
```

then `ecpg` will already do the substitution and your C compiler will never see any name or identifier `MYNUMBER`. Note that you cannot use `#define` for a constant that you are going to use in an embedded SQL query because in this case the embedded SQL precompiler is not able to see this declaration.

If multiple input files are named on the `ecpg` preprocessor's command line, the effects of `EXEC SQL DEFINE` and `EXEC SQL UNDEF` do not carry across files: each file starts with only the symbols defined by `-D` switches on the command line.

37.9.3. `ifdef`, `ifndef`, `elif`, `else`, and `endif` Directives

You can use the following directives to compile code sections conditionally:

```
EXEC SQL ifdef name;
```

Checks a *name* and processes subsequent lines if *name* has been defined via `EXEC SQL define name`.

```
EXEC SQL ifndef name;
```

Checks a *name* and processes subsequent lines if *name* has *not* been defined via `EXEC SQL define name`.

```
EXEC SQL elif name;
```

Begins an optional alternative section after an `EXEC SQL ifdef name` or `EXEC SQL ifndef name` directive. Any number of `elif` sections can appear. Lines following an `elif` will be processed if

name has been defined *and* no previous section of the same `ifdef/ifndef...endif` construct has been processed.

```
EXEC SQL else;
```

Begins an optional, final alternative section after an `EXEC SQL ifdef name` or `EXEC SQL ifndef name` directive. Subsequent lines will be processed if no previous section of the same `ifdef/ifndef...endif` construct has been processed.

```
EXEC SQL endif;
```

Ends an `ifdef/ifndef...endif` construct. Subsequent lines are processed normally.

`ifdef/ifndef...endif` constructs can be nested, up to 127 levels deep.

This example will compile exactly one of the three `SET TIMEZONE` commands:

```
EXEC SQL ifdef TZVAR;
EXEC SQL SET TIMEZONE TO TZVAR;
EXEC SQL elif TZNAME;
EXEC SQL SET TIMEZONE TO TZNAME;
EXEC SQL else;
EXEC SQL SET TIMEZONE TO 'GMT';
EXEC SQL endif;
```

37.10. Processing Embedded SQL Programs

Now that you have an idea how to form embedded SQL C programs, you probably want to know how to compile them. Before compiling you run the file through the embedded SQL C preprocessor, which converts the SQL statements you used to special function calls. After compiling, you must link with a special library that contains the needed functions. These functions fetch information from the arguments, perform the SQL command using the `libpq` interface, and put the result in the arguments specified for output.

The preprocessor program is called `ecpg` and is included in a normal Postgres Pro installation. Embedded SQL programs are typically named with an extension `.pgc`. If you have a program file called `prog1.pgc`, you can preprocess it by simply calling:

```
ecpg prog1.pgc
```

This will create a file called `prog1.c`. If your input files do not follow the suggested naming pattern, you can specify the output file explicitly using the `-o` option.

The preprocessed file can be compiled normally, for example:

```
cc -c prog1.c
```

The generated C source files include header files from the Postgres Pro installation, so if you installed Postgres Pro in a location that is not searched by default, you have to add an option such as `-I/usr/local/pgsql/include` to the compilation command line.

To link an embedded SQL program, you need to include the `libecpg` library, like so:

```
cc -o myprog prog1.o prog2.o ... -lecpg
```

Again, you might have to add an option like `-L/usr/local/pgsql/lib` to that command line.

You can use `pg_config` or `pkg-config` with package name `libecpg` to get the paths for your installation.

If you manage the build process of a larger project using `make`, it might be convenient to include the following implicit rule to your makefiles:

```
ECPG = ecpg
```

```
%c: %.pgc
$(ECPG) $<
```

The complete syntax of the `ecpg` command is detailed in [ecpg](#).

The `ecpg` library is thread-safe by default. However, you might need to use some threading command-line options to compile your client code.

37.11. Library Functions

The `libecpg` library primarily contains “hidden” functions that are used to implement the functionality expressed by the embedded SQL commands. But there are some functions that can usefully be called directly. Note that this makes your code unportable.

- `ECPGdebug(int on, FILE *stream)` turns on debug logging if called with the first argument non-zero. Debug logging is done on *stream*. The log contains all SQL statements with all the input variables inserted, and the results from the Postgres Pro server. This can be very useful when searching for errors in your SQL statements.

Note

On Windows, if the `ecpg` libraries and an application are compiled with different flags, this function call will crash the application because the internal representation of the `FILE` pointers differ. Specifically, multithreaded/single-threaded, release/debug, and static/dynamic flags should be the same for the library and all applications using that library.

- `ECPGget_PGconn(const char *connection_name)` returns the library database connection handle identified by the given name. If *connection_name* is set to `NULL`, the current connection handle is returned. If no connection handle can be identified, the function returns `NULL`. The returned connection handle can be used to call any other functions from `libpq`, if necessary.

Note

It is a bad idea to manipulate database connection handles made from `ecpg` directly with `libpq` routines.

- `ECPGtransactionStatus(const char *connection_name)` returns the current transaction status of the given connection identified by *connection_name*. See [Section 35.2](#) and `libpq`'s `PQtransactionStatus` for details about the returned status codes.
- `ECPGstatus(int lineno, const char* connection_name)` returns true if you are connected to a database and false if not. *connection_name* can be `NULL` if a single connection is being used.

37.12. Large Objects

Large objects are not directly supported by ECPG, but ECPG application can manipulate large objects through the `libpq` large object functions, obtaining the necessary `PGconn` object by calling the `ECPGget_PGconn()` function. (However, use of the `ECPGget_PGconn()` function and touching `PGconn` objects directly should be done very carefully and ideally not mixed with other ECPG database access calls.)

For more details about the `ECPGget_PGconn()`, see [Section 37.11](#). For information about the large object function interface, see [Chapter 36](#).

Large object functions have to be called in a transaction block, so when `autocommit` is off, `BEGIN` commands have to be issued explicitly.

[Example 37.2](#) shows an example program that illustrates how to create, write, and read a large object in an ECPG application.

Example 37.2. ECPG Program Accessing Large Objects

```
#include <stdio.h>
#include <stdlib.h>
#include <libpq-fe.h>
#include <libpq/libpq-fs.h>

EXEC SQL WHENEVER SQLERROR STOP;

int
main(void)
{
    PGconn      *conn;
    Oid          loid;
    int          fd;
    char         buf[256];
    int          buflen = 256;
    char         buf2[256];
    int          rc;

    memset(buf, 1, buflen);

    EXEC SQL CONNECT TO testdb AS con1;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    conn = ECPGget_PGconn("con1");
    printf("conn = %p\n", conn);

    /* create */
    loid = lo_create(conn, 0);
    if (loid < 0)
        printf("lo_create() failed: %s", PQerrorMessage(conn));

    printf("loid = %d\n", loid);

    /* write test */
    fd = lo_open(conn, loid, INV_READ|INV_WRITE);
    if (fd < 0)
        printf("lo_open() failed: %s", PQerrorMessage(conn));

    printf("fd = %d\n", fd);

    rc = lo_write(conn, fd, buf, buflen);
    if (rc < 0)
        printf("lo_write() failed\n");

    rc = lo_close(conn, fd);
    if (rc < 0)
        printf("lo_close() failed: %s", PQerrorMessage(conn));

    /* read test */
    fd = lo_open(conn, loid, INV_READ);
    if (fd < 0)
        printf("lo_open() failed: %s", PQerrorMessage(conn));

    printf("fd = %d\n", fd);
```

```
rc = lo_read(conn, fd, buf2, buflen);
if (rc < 0)
    printf("lo_read() failed\n");

rc = lo_close(conn, fd);
if (rc < 0)
    printf("lo_close() failed: %s", PQerrorMessage(conn));

/* check */
rc = memcmp(buf, buf2, buflen);
printf("memcmp() = %d\n", rc);

/* cleanup */
rc = lo_unlink(conn, loid);
if (rc < 0)
    printf("lo_unlink() failed: %s", PQerrorMessage(conn));

EXEC SQL COMMIT;
EXEC SQL DISCONNECT ALL;
return 0;
}
```

37.13. C++ Applications

ECPG has some limited support for C++ applications. This section describes some caveats.

The `ecpg` preprocessor takes an input file written in C (or something like C) and embedded SQL commands, converts the embedded SQL commands into C language chunks, and finally generates a `.c` file. The header file declarations of the library functions used by the C language chunks that `ecpg` generates are wrapped in `extern "C" { ... }` blocks when used under C++, so they should work seamlessly in C++.

In general, however, the `ecpg` preprocessor only understands C; it does not handle the special syntax and reserved words of the C++ language. So, some embedded SQL code written in C++ application code that uses complicated features specific to C++ might fail to be preprocessed correctly or might not work as expected.

A safe way to use the embedded SQL code in a C++ application is hiding the ECPG calls in a C module, which the C++ application code calls into to access the database, and linking that together with the rest of the C++ code. See [Section 37.13.2](#) about that.

37.13.1. Scope for Host Variables

The `ecpg` preprocessor understands the scope of variables in C. In the C language, this is rather simple because the scopes of variables is based on their code blocks. In C++, however, the class member variables are referenced in a different code block from the declared position, so the `ecpg` preprocessor will not understand the scope of the class member variables.

For example, in the following case, the `ecpg` preprocessor cannot find any declaration for the variable `dbname` in the `test` method, so an error will occur.

```
class TestCpp
{
    EXEC SQL BEGIN DECLARE SECTION;
    char dbname[1024];
    EXEC SQL END DECLARE SECTION;

public:
```

```
TestCpp();
void test();
~TestCpp();
};

TestCpp::TestCpp()
{
    EXEC SQL CONNECT TO testdb1;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
}

void Test::test()
{
    EXEC SQL SELECT current_database() INTO :dbname;
    printf("current_database = %s\n", dbname);
}

TestCpp::~TestCpp()
{
    EXEC SQL DISCONNECT ALL;
}
```

This code will result in an error like this:

ecpg test_cpp.pgc

test_cpp.pgc:28: ERROR: variable "dbname" is not declared

To avoid this scope issue, the `test` method could be modified to use a local variable as intermediate storage. But this approach is only a poor workaround, because it uglifies the code and reduces performance.

```
void TestCpp::test()
{
    EXEC SQL BEGIN DECLARE SECTION;
    char tmp[1024];
    EXEC SQL END DECLARE SECTION;

    EXEC SQL SELECT current_database() INTO :tmp;
    strcpy(dbname, tmp, sizeof(tmp));

    printf("current_database = %s\n", dbname);
}
```

37.13.2. C++ Application Development with External C Module

If you understand these technical limitations of the `ecpg` preprocessor in C++, you might come to the conclusion that linking C objects and C++ objects at the link stage to enable C++ applications to use ECPG features could be better than writing some embedded SQL commands in C++ code directly. This section describes a way to separate some embedded SQL commands from C++ application code with a simple example. In this example, the application is implemented in C++, while C and ECPG is used to connect to the Postgres Pro server.

Three kinds of files have to be created: a C file (*.pgc), a header file, and a C++ file:

test_mod.pgc

A sub-routine module to execute SQL commands embedded in C. It is going to be converted into `test_mod.c` by the preprocessor.

```
#include "test_mod.h"
#include <stdio.h>

void
```



```
db_connect ()
{
    EXEC SQL CONNECT TO testdb1;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL
    COMMIT;
}

void
db_test ()
{
    EXEC SQL BEGIN DECLARE SECTION;
    char dbname[1024];
    EXEC SQL END DECLARE SECTION;

    EXEC SQL SELECT current_database() INTO :dbname;
    printf("current_database = %s\n", dbname);
}

void
db_disconnect ()
{
    EXEC SQL DISCONNECT ALL;
}
```

test_mod.h

A header file with declarations of the functions in the C module (`test_mod.pgc`). It is included by `test_cpp.cpp`. This file has to have an `extern "C"` block around the declarations, because it will be linked from the C++ module.

```
#ifdef __cplusplus
extern "C" {
#endif

void db_connect ();
void db_test ();
void db_disconnect ();

#ifdef __cplusplus
}
#endif
```

test_cpp.cpp

The main code for the application, including the `main` routine, and in this example a C++ class.

```
#include "test_mod.h"

class TestCpp
{
public:
    TestCpp();
    void test();
    ~TestCpp();
};

TestCpp::TestCpp()
{
    db_connect();
}
```

```
void
TestCpp::test()
{
    db_test();
}

TestCpp::~TestCpp()
{
    db_disconnect();
}

int
main(void)
{
    TestCpp *t = new TestCpp();

    t->test();
    return 0;
}
```

To build the application, proceed as follows. Convert `test_mod.pgc` into `test_mod.c` by running `ecpg`, and generate `test_mod.o` by compiling `test_mod.c` with the C compiler:

```
ecpg -o test_mod.c test_mod.pgc
cc -c test_mod.c -o test_mod.o
```

Next, generate `test_cpp.o` by compiling `test_cpp.cpp` with the C++ compiler:

```
c++ -c test_cpp.cpp -o test_cpp.o
```

Finally, link these object files, `test_cpp.o` and `test_mod.o`, into one executable, using the C++ compiler driver:

```
c++ test_cpp.o test_mod.o -lecp -o test_cpp
```

37.14. Embedded SQL Commands

This section describes all SQL commands that are specific to embedded SQL. Also refer to the SQL commands listed in [SQL Commands](#), which can also be used in embedded SQL, unless stated otherwise.

ALLOCATE DESCRIPTOR

ALLOCATE DESCRIPTOR — allocate an SQL descriptor area

Synopsis

```
ALLOCATE DESCRIPTOR name
```

Description

ALLOCATE DESCRIPTOR allocates a new named SQL descriptor area, which can be used to exchange data between the Postgres Pro server and the host program.

Descriptor areas should be freed after use using the DEALLOCATE DESCRIPTOR command.

Parameters

name

A name of SQL descriptor, case sensitive. This can be an SQL identifier or a host variable.

Examples

```
EXEC SQL ALLOCATE DESCRIPTOR mydesc;
```

Compatibility

ALLOCATE DESCRIPTOR is specified in the SQL standard.

See Also

[DEALLOCATE DESCRIPTOR](#), [GET DESCRIPTOR](#), [SET DESCRIPTOR](#)

CONNECT

CONNECT — establish a database connection

Synopsis

```
CONNECT TO connection_target [ AS connection_name ] [ USER connection_user ]  
CONNECT TO DEFAULT  
CONNECT connection_user  
DATABASE connection_target
```

Description

The `CONNECT` command establishes a connection between the client and the Postgres Pro server.

Parameters

connection_target

connection_target specifies the target server of the connection on one of several forms.

[*database_name*] [@*host*] [:*port*]

Connect over TCP/IP

unix:postgresql://*host* [:*port*] / [*database_name*] [?*connection_option*]

Connect over Unix-domain sockets

tcp:postgresql://*host* [:*port*] / [*database_name*] [?*connection_option*]

Connect over TCP/IP

SQL string constant

containing a value in one of the above forms

host variable

host variable of type `char[]` or `VARCHAR[]` containing a value in one of the above forms

connection_name

An optional identifier for the connection, so that it can be referred to in other commands. This can be an SQL identifier or a host variable.

connection_user

The user name for the database connection.

This parameter can also specify user name and password, using one the forms *user_name/password*, *user_name IDENTIFIED BY password*, or *user_name USING password*.

User name and password can be SQL identifiers, string constants, or host variables.

DEFAULT

Use all default connection parameters, as defined by libpq.

Examples

Here a several variants for specifying connection parameters:

```
EXEC SQL CONNECT TO "connectdb" AS main;  
EXEC SQL CONNECT TO "connectdb" AS second;
```

```
EXEC SQL CONNECT TO "unix:postgresql://200.46.204.71/connectdb" AS main USER
connectuser;
EXEC SQL CONNECT TO "unix:postgresql://localhost/connectdb" AS main USER connectuser;
EXEC SQL CONNECT TO 'connectdb' AS main;
EXEC SQL CONNECT TO 'unix:postgresql://localhost/connectdb' AS main USER :user;
EXEC SQL CONNECT TO :db AS :id;
EXEC SQL CONNECT TO :db USER connectuser USING :pw;
EXEC SQL CONNECT TO @localhost AS main USER connectdb;
EXEC SQL CONNECT TO REGRESSDB1 as main;
EXEC SQL CONNECT TO AS main USER connectdb;
EXEC SQL CONNECT TO connectdb AS :id;
EXEC SQL CONNECT TO connectdb AS main USER connectuser/connectdb;
EXEC SQL CONNECT TO connectdb AS main;
EXEC SQL CONNECT TO connectdb@localhost AS main;
EXEC SQL CONNECT TO tcp:postgresql://localhost/ USER connectdb;
EXEC SQL CONNECT TO tcp:postgresql://localhost/connectdb USER connectuser IDENTIFIED BY
connectpw;
EXEC SQL CONNECT TO tcp:postgresql://localhost:20/connectdb USER connectuser IDENTIFIED
BY connectpw;
EXEC SQL CONNECT TO unix:postgresql://localhost/ AS main USER connectdb;
EXEC SQL CONNECT TO unix:postgresql://localhost/connectdb AS main USER connectuser;
EXEC SQL CONNECT TO unix:postgresql://localhost/connectdb USER connectuser IDENTIFIED
BY "connectpw";
EXEC SQL CONNECT TO unix:postgresql://localhost/connectdb USER connectuser USING
"connectpw";
EXEC SQL CONNECT TO unix:postgresql://localhost/connectdb?connect_timeout=14 USER
connectuser;
```

Here is an example program that illustrates the use of host variables to specify connection parameters:

```
int
main(void)
{
EXEC SQL BEGIN DECLARE SECTION;
    char *dbname      = "testdb";      /* database name */
    char *user        = "testuser";    /* connection user name */
    char *connection  = "tcp:postgresql://localhost:5432/testdb";
                                      /* connection string */
    char ver[256];      /* buffer to store the version string */
EXEC SQL END DECLARE SECTION;

    ECPGdebug(1, stderr);

    EXEC SQL CONNECT TO :dbname USER :user;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
    EXEC SQL SELECT pgpro_version() INTO :ver;
    EXEC SQL DISCONNECT;

    printf("version: %s\n", ver);

    EXEC SQL CONNECT TO :connection USER :user;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
    EXEC SQL SELECT pgpro_version() INTO :ver;
    EXEC SQL DISCONNECT;

    printf("version: %s\n", ver);

    return 0;
}
```

Compatibility

`CONNECT` is specified in the SQL standard, but the format of the connection parameters is implementation-specific.

See Also

[DISCONNECT](#), [SET CONNECTION](#)

DEALLOCATE DESCRIPTOR

DEALLOCATE DESCRIPTOR — deallocate an SQL descriptor area

Synopsis

```
DEALLOCATE DESCRIPTOR name
```

Description

DEALLOCATE DESCRIPTOR deallocates a named SQL descriptor area.

Parameters

name

The name of the descriptor which is going to be deallocated. It is case sensitive. This can be an SQL identifier or a host variable.

Examples

```
EXEC SQL DEALLOCATE DESCRIPTOR mydesc;
```

Compatibility

DEALLOCATE DESCRIPTOR is specified in the SQL standard.

See Also

[ALLOCATE DESCRIPTOR](#), [GET DESCRIPTOR](#), [SET DESCRIPTOR](#)

DECLARE

DECLARE — define a cursor

Synopsis

```
DECLARE cursor_name [ BINARY ] [ ASENSITIVE | INSENSITIVE ] [ [ NO ] SCROLL ] CURSOR  
  [ { WITH | WITHOUT } HOLD ] FOR prepared_name  
DECLARE cursor_name [ BINARY ] [ ASENSITIVE | INSENSITIVE ] [ [ NO ] SCROLL ] CURSOR  
  [ { WITH | WITHOUT } HOLD ] FOR query
```

Description

DECLARE declares a cursor for iterating over the result set of a prepared statement. This command has slightly different semantics from the direct SQL command `DECLARE`: Whereas the latter executes a query and prepares the result set for retrieval, this embedded SQL command merely declares a name as a “loop variable” for iterating over the result set of a query; the actual execution happens when the cursor is opened with the `OPEN` command.

Parameters

cursor_name

A cursor name, case sensitive. This can be an SQL identifier or a host variable.

prepared_name

The name of a prepared query, either as an SQL identifier or a host variable.

query

A [SELECT](#) or [VALUES](#) command which will provide the rows to be returned by the cursor.

For the meaning of the cursor options, see [DECLARE](#).

Examples

Examples declaring a cursor for a query:

```
EXEC SQL DECLARE C CURSOR FOR SELECT * FROM My_Table;  
EXEC SQL DECLARE C CURSOR FOR SELECT Item1 FROM T;  
EXEC SQL DECLARE curl CURSOR FOR SELECT pgpro_version();
```

An example declaring a cursor for a prepared statement:

```
EXEC SQL PREPARE stmt1 AS SELECT pgpro_version();  
EXEC SQL DECLARE curl CURSOR FOR stmt1;
```

Compatibility

DECLARE is specified in the SQL standard.

See Also

[OPEN](#), [CLOSE](#), [DECLARE](#)

DECLARE STATEMENT

DECLARE STATEMENT — declare SQL statement identifier

Synopsis

```
EXEC SQL [ AT connection_name ] DECLARE statement_name STATEMENT
```

Description

DECLARE STATEMENT declares an SQL statement identifier. SQL statement identifier can be associated with the connection. When the identifier is used by dynamic SQL statements, the statements are executed using the associated connection. The namespace of the declaration is the precompile unit, and multiple declarations to the same SQL statement identifier are not allowed. Note that if the precompiler runs in Informix compatibility mode and some SQL statement is declared, "database" can not be used as a cursor name.

Parameters

connection_name

A database connection name established by the CONNECT command.

AT clause can be omitted, but such statement has no meaning.

statement_name

The name of an SQL statement identifier, either as an SQL identifier or a host variable.

Notes

This association is valid only if the declaration is physically placed on top of a dynamic statement.

Examples

```
EXEC SQL CONNECT TO postgres AS con1;
EXEC SQL AT con1 DECLARE sql_stmt STATEMENT;
EXEC SQL DECLARE cursor_name CURSOR FOR sql_stmt;
EXEC SQL PREPARE sql_stmt FROM :dyn_string;
EXEC SQL OPEN cursor_name;
EXEC SQL FETCH cursor_name INTO :column1;
EXEC SQL CLOSE cursor_name;
```

Compatibility

DECLARE STATEMENT is an extension of the SQL standard, but can be used in famous DBMSs.

See Also

[CONNECT](#), [DECLARE](#), [OPEN](#)

DESCRIBE

DESCRIBE — obtain information about a prepared statement or result set

Synopsis

```
DESCRIBE [ OUTPUT ] prepared_name USING [ SQL ] DESCRIPTOR descriptor_name
DESCRIBE [ OUTPUT ] prepared_name INTO [ SQL ] DESCRIPTOR descriptor_name
DESCRIBE [ OUTPUT ] prepared_name INTO sqlda_name
```

Description

DESCRIBE retrieves metadata information about the result columns contained in a prepared statement, without actually fetching a row.

Parameters

prepared_name

The name of a prepared statement. This can be an SQL identifier or a host variable.

descriptor_name

A descriptor name. It is case sensitive. It can be an SQL identifier or a host variable.

sqlda_name

The name of an SQLDA variable.

Examples

```
EXEC SQL ALLOCATE DESCRIPTOR mydesc;
EXEC SQL PREPARE stmt1 FROM :sql_stmt;
EXEC SQL DESCRIBE stmt1 INTO SQL DESCRIPTOR mydesc;
EXEC SQL GET DESCRIPTOR mydesc VALUE 1 :charvar = NAME;
EXEC SQL DEALLOCATE DESCRIPTOR mydesc;
```

Compatibility

DESCRIBE is specified in the SQL standard.

See Also

[ALLOCATE DESCRIPTOR](#), [GET DESCRIPTOR](#)

DISCONNECT

DISCONNECT — terminate a database connection

Synopsis

```
DISCONNECT connection_name
DISCONNECT [ CURRENT ]
DISCONNECT ALL
```

Description

DISCONNECT closes a connection (or all connections) to the database.

Parameters

connection_name

A database connection name established by the `CONNECT` command.

CURRENT

Close the “current” connection, which is either the most recently opened connection, or the connection set by the `SET CONNECTION` command. This is also the default if no argument is given to the `DISCONNECT` command.

ALL

Close all open connections.

Examples

```
int
main(void)
{
    EXEC SQL CONNECT TO testdb AS con1 USER testuser;
    EXEC SQL CONNECT TO testdb AS con2 USER testuser;
    EXEC SQL CONNECT TO testdb AS con3 USER testuser;

    EXEC SQL DISCONNECT CURRENT; /* close con3          */
    EXEC SQL DISCONNECT ALL;      /* close con2 and con1 */

    return 0;
}
```

Compatibility

DISCONNECT is specified in the SQL standard.

See Also

[CONNECT](#), [SET CONNECTION](#)

EXECUTE IMMEDIATE

EXECUTE IMMEDIATE — dynamically prepare and execute a statement

Synopsis

```
EXECUTE IMMEDIATE string
```

Description

EXECUTE IMMEDIATE immediately prepares and executes a dynamically specified SQL statement, without retrieving result rows.

Parameters

string

A literal string or a host variable containing the SQL statement to be executed.

Notes

In typical usage, the *string* is a host variable reference to a string containing a dynamically-constructed SQL statement. The case of a literal string is not very useful; you might as well just write the SQL statement directly, without the extra typing of EXECUTE IMMEDIATE.

If you do use a literal string, keep in mind that any double quotes you might wish to include in the SQL statement must be written as octal escapes (`\042`) not the usual C idiom `\"`. This is because the string is inside an EXEC SQL section, so the ECPG lexer parses it according to SQL rules not C rules. Any embedded backslashes will later be handled according to C rules; but `\"` causes an immediate syntax error because it is seen as ending the literal.

Examples

Here is an example that executes an INSERT statement using EXECUTE IMMEDIATE and a host variable named `command`:

```
sprintf(command, "INSERT INTO test (name, amount, letter) VALUES ('db: 'r1'', 1, 'f')");  
EXEC SQL EXECUTE IMMEDIATE :command;
```

Compatibility

EXECUTE IMMEDIATE is specified in the SQL standard.

GET DESCRIPTOR

GET DESCRIPTOR — get information from an SQL descriptor area

Synopsis

```
GET DESCRIPTOR descriptor_name :cvariable = descriptor_header_item [, ... ]
GET DESCRIPTOR descriptor_name VALUE column_number :cvariable = descriptor_item
[, ... ]
```

Description

GET DESCRIPTOR retrieves information about a query result set from an SQL descriptor area and stores it into host variables. A descriptor area is typically populated using `FETCH` or `SELECT` before using this command to transfer the information into host language variables.

This command has two forms: The first form retrieves descriptor “header” items, which apply to the result set in its entirety. One example is the row count. The second form, which requires the column number as additional parameter, retrieves information about a particular column. Examples are the column name and the actual column value.

Parameters

descriptor_name

A descriptor name.

descriptor_header_item

A token identifying which header information item to retrieve. Only `COUNT`, to get the number of columns in the result set, is currently supported.

column_number

The number of the column about which information is to be retrieved. The count starts at 1.

descriptor_item

A token identifying which item of information about a column to retrieve. See [Section 37.7.1](#) for a list of supported items.

cvariable

A host variable that will receive the data retrieved from the descriptor area.

Examples

An example to retrieve the number of columns in a result set:

```
EXEC SQL GET DESCRIPTOR d :d_count = COUNT;
```

An example to retrieve a data length in the first column:

```
EXEC SQL GET DESCRIPTOR d VALUE 1 :d_returned_octet_length = RETURNED_OCTET_LENGTH;
```

An example to retrieve the data body of the second column as a string:

```
EXEC SQL GET DESCRIPTOR d VALUE 2 :d_data = DATA;
```

Here is an example for a whole procedure of executing `SELECT current_database();` and showing the number of columns, the column data length, and the column data:

```
int
main(void)
{
```

```
EXEC SQL BEGIN DECLARE SECTION;
    int  d_count;
    char d_data[1024];
    int  d_returned_octet_length;
EXEC SQL END DECLARE SECTION;

EXEC SQL CONNECT TO testdb AS con1 USER testuser;
EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
EXEC SQL ALLOCATE DESCRIPTOR d;

/* Declare, open a cursor, and assign a descriptor to the cursor */
EXEC SQL DECLARE cur CURSOR FOR SELECT current_database();
EXEC SQL OPEN cur;
EXEC SQL FETCH NEXT FROM cur INTO SQL DESCRIPTOR d;

/* Get a number of total columns */
EXEC SQL GET DESCRIPTOR d :d_count = COUNT;
printf("d_count          = %d\n", d_count);

/* Get length of a returned column */
EXEC SQL GET DESCRIPTOR d VALUE 1 :d_returned_octet_length = RETURNED_OCTET_LENGTH;
printf("d_returned_octet_length = %d\n", d_returned_octet_length);

/* Fetch the returned column as a string */
EXEC SQL GET DESCRIPTOR d VALUE 1 :d_data = DATA;
printf("d_data          = %s\n", d_data);

/* Closing */
EXEC SQL CLOSE cur;
EXEC SQL COMMIT;

EXEC SQL DEALLOCATE DESCRIPTOR d;
EXEC SQL DISCONNECT ALL;

return 0;
}
```

When the example is executed, the result will look like this:

```
d_count          = 1
d_returned_octet_length = 6
d_data          = testdb
```

Compatibility

GET DESCRIPTOR is specified in the SQL standard.

See Also

[ALLOCATE DESCRIPTOR](#), [SET DESCRIPTOR](#)

OPEN

OPEN — open a dynamic cursor

Synopsis

```
OPEN cursor_name
OPEN cursor_name USING value [, ... ]
OPEN cursor_name USING SQL DESCRIPTOR descriptor_name
```

Description

OPEN opens a cursor and optionally binds actual values to the placeholders in the cursor's declaration. The cursor must previously have been declared with the `DECLARE` command. The execution of OPEN causes the query to start executing on the server.

Parameters

cursor_name

The name of the cursor to be opened. This can be an SQL identifier or a host variable.

value

A value to be bound to a placeholder in the cursor. This can be an SQL constant, a host variable, or a host variable with indicator.

descriptor_name

The name of a descriptor containing values to be bound to the placeholders in the cursor. This can be an SQL identifier or a host variable.

Examples

```
EXEC SQL OPEN a;
EXEC SQL OPEN d USING 1, 'test';
EXEC SQL OPEN c1 USING SQL DESCRIPTOR mydesc;
EXEC SQL OPEN :curname1;
```

Compatibility

OPEN is specified in the SQL standard.

See Also

[DECLARE](#), [CLOSE](#)

PREPARE

PREPARE — prepare a statement for execution

Synopsis

```
PREPARE prepared_name FROM string
```

Description

PREPARE prepares a statement dynamically specified as a string for execution. This is different from the direct SQL statement [PREPARE](#), which can also be used in embedded programs. The [EXECUTE](#) command is used to execute either kind of prepared statement.

Parameters

prepared_name

An identifier for the prepared query.

string

A literal string or a host variable containing a preparable SQL statement, one of SELECT, INSERT, UPDATE, or DELETE. Use question marks (?) for parameter values to be supplied at execution.

Notes

In typical usage, the *string* is a host variable reference to a string containing a dynamically-constructed SQL statement. The case of a literal string is not very useful; you might as well just write a direct SQL PREPARE statement.

If you do use a literal string, keep in mind that any double quotes you might wish to include in the SQL statement must be written as octal escapes (\042) not the usual C idiom \". This is because the string is inside an EXEC SQL section, so the ECPG lexer parses it according to SQL rules not C rules. Any embedded backslashes will later be handled according to C rules; but \" causes an immediate syntax error because it is seen as ending the literal.

Examples

```
char *stmt = "SELECT * FROM test1 WHERE a = ? AND b = ?";
```

```
EXEC SQL ALLOCATE DESCRIPTOR outdesc;  
EXEC SQL PREPARE foo FROM :stmt;
```

```
EXEC SQL EXECUTE foo USING SQL DESCRIPTOR indesc INTO SQL DESCRIPTOR outdesc;
```

Compatibility

PREPARE is specified in the SQL standard.

See Also

[EXECUTE](#)

SET AUTOCOMMIT

SET AUTOCOMMIT — set the autocommit behavior of the current session

Synopsis

```
SET AUTOCOMMIT { = | TO } { ON | OFF }
```

Description

SET AUTOCOMMIT sets the autocommit behavior of the current database session. By default, embedded SQL programs are *not* in autocommit mode, so COMMIT needs to be issued explicitly when desired. This command can change the session to autocommit mode, where each individual statement is committed implicitly.

Compatibility

SET AUTOCOMMIT is an extension of Postgres Pro ECPG.

SET CONNECTION

SET CONNECTION — select a database connection

Synopsis

```
SET CONNECTION [ TO | = ] connection_name
```

Description

SET CONNECTION sets the “current” database connection, which is the one that all commands use unless overridden.

Parameters

connection_name

A database connection name established by the `CONNECT` command.

CURRENT

Set the connection to the current connection (thus, nothing happens).

Examples

```
EXEC SQL SET CONNECTION TO con2;  
EXEC SQL SET CONNECTION = con1;
```

Compatibility

SET CONNECTION is specified in the SQL standard.

See Also

[CONNECT](#), [DISCONNECT](#)

SET DESCRIPTOR

SET DESCRIPTOR — set information in an SQL descriptor area

Synopsis

```
SET DESCRIPTOR descriptor_name descriptor_header_item = value [, ... ]  
SET DESCRIPTOR descriptor_name VALUE number descriptor_item = value [, ...]
```

Description

SET DESCRIPTOR populates an SQL descriptor area with values. The descriptor area is then typically used to bind parameters in a prepared query execution.

This command has two forms: The first form applies to the descriptor “header”, which is independent of a particular datum. The second form assigns values to particular datums, identified by number.

Parameters

descriptor_name

A descriptor name.

descriptor_header_item

A token identifying which header information item to set. Only COUNT, to set the number of descriptor items, is currently supported.

number

The number of the descriptor item to set. The count starts at 1.

descriptor_item

A token identifying which item of information to set in the descriptor. See [Section 37.7.1](#) for a list of supported items.

value

A value to store into the descriptor item. This can be an SQL constant or a host variable.

Examples

```
EXEC SQL SET DESCRIPTOR indesc COUNT = 1;  
EXEC SQL SET DESCRIPTOR indesc VALUE 1 DATA = 2;  
EXEC SQL SET DESCRIPTOR indesc VALUE 1 DATA = :val1;  
EXEC SQL SET DESCRIPTOR indesc VALUE 2 INDICATOR = :val1, DATA = 'some string';  
EXEC SQL SET DESCRIPTOR indesc VALUE 2 INDICATOR = :val2null, DATA = :val2;
```

Compatibility

SET DESCRIPTOR is specified in the SQL standard.

See Also

[ALLOCATE DESCRIPTOR](#), [GET DESCRIPTOR](#)

TYPE

TYPE — define a new data type

Synopsis

```
TYPE type_name IS ctype
```

Description

The TYPE command defines a new C type. It is equivalent to putting a `typedef` into a declare section.

This command is only recognized when `ecpg` is run with the `-c` option.

Parameters

type_name

The name for the new type. It must be a valid C type name.

ctype

A C type specification.

Examples

```
EXEC SQL TYPE customer IS
    struct
    {
        varchar name[50];
        int      phone;
    };

```

```
EXEC SQL TYPE cust_ind IS
    struct ind
    {
        short  name_ind;
        short  phone_ind;
    };

```

```
EXEC SQL TYPE c IS char reference;
EXEC SQL TYPE ind IS union { int integer; short smallint; };
EXEC SQL TYPE intarray IS int[AMOUNT];
EXEC SQL TYPE str IS varchar[BUFFERSIZ];
EXEC SQL TYPE string IS char[11];

```

Here is an example program that uses EXEC SQL TYPE:

```
EXEC SQL WHENEVER SQLERROR SQLPRINT;

EXEC SQL TYPE tt IS
    struct
    {
        varchar v[256];
        int      i;
    };

EXEC SQL TYPE tt_ind IS
    struct ind {
        short  v_ind;
        short  i_ind;
    };

```

```
int
main(void)
{
EXEC SQL BEGIN DECLARE SECTION;
    tt t;
    tt_ind t_ind;
EXEC SQL END DECLARE SECTION;

    EXEC SQL CONNECT TO testdb AS con1;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;

    EXEC SQL SELECT current_database(), 256 INTO :t:t_ind LIMIT 1;

    printf("t.v = %s\n", t.v.arr);
    printf("t.i = %d\n", t.i);

    printf("t_ind.v_ind = %d\n", t_ind.v_ind);
    printf("t_ind.i_ind = %d\n", t_ind.i_ind);

    EXEC SQL DISCONNECT con1;

    return 0;
}
```

The output from this program looks like this:

```
t.v = testdb
t.i = 256
t_ind.v_ind = 0
t_ind.i_ind = 0
```

Compatibility

The `TYPE` command is a Postgres Pro extension.

VAR

VAR — define a variable

Synopsis

```
VAR varname IS ctype
```

Description

The VAR command assigns a new C data type to a host variable. The host variable must be previously declared in a declare section.

Parameters

varname

A C variable name.

ctype

A C type specification.

Examples

```
Exec sql begin declare section;  
short a;  
exec sql end declare section;  
EXEC SQL VAR a IS int;
```

Compatibility

The VAR command is a Postgres Pro extension.

WHENEVER

WHENEVER — specify the action to be taken when an SQL statement causes a specific class condition to be raised

Synopsis

```
WHENEVER { NOT FOUND | SQLERROR | SQLWARNING } action
```

Description

Define a behavior which is called on the special cases (Rows not found, SQL warnings or errors) in the result of SQL execution.

Parameters

See [Section 37.8.1](#) for a description of the parameters.

Examples

```
EXEC SQL WHENEVER NOT FOUND CONTINUE;
EXEC SQL WHENEVER NOT FOUND DO BREAK;
EXEC SQL WHENEVER NOT FOUND DO CONTINUE;
EXEC SQL WHENEVER SQLWARNING SQLPRINT;
EXEC SQL WHENEVER SQLWARNING DO warn();
EXEC SQL WHENEVER SQLERROR sqlprint;
EXEC SQL WHENEVER SQLERROR CALL print2();
EXEC SQL WHENEVER SQLERROR DO handle_error("select");
EXEC SQL WHENEVER SQLERROR DO sqlnotice(NULL, NONO);
EXEC SQL WHENEVER SQLERROR DO sqlprint();
EXEC SQL WHENEVER SQLERROR GOTO error_label;
EXEC SQL WHENEVER SQLERROR STOP;
```

A typical application is the use of WHENEVER NOT FOUND BREAK to handle looping through result sets:

```
int
main(void)
{
    EXEC SQL CONNECT TO testdb AS con1;
    EXEC SQL SELECT pg_catalog.set_config('search_path', '', false); EXEC SQL COMMIT;
    EXEC SQL ALLOCATE DESCRIPTOR d;
    EXEC SQL DECLARE cur CURSOR FOR SELECT current_database(), 'hoge', 256;
    EXEC SQL OPEN cur;

    /* when end of result set reached, break out of while loop */
    EXEC SQL WHENEVER NOT FOUND DO BREAK;

    while (1)
    {
        EXEC SQL FETCH NEXT FROM cur INTO SQL DESCRIPTOR d;
        ...
    }

    EXEC SQL CLOSE cur;
    EXEC SQL COMMIT;

    EXEC SQL DEALLOCATE DESCRIPTOR d;
    EXEC SQL DISCONNECT ALL;

    return 0;
}
```

Compatibility

WHENEVER is specified in the SQL standard, but most of the actions are Postgres Pro extensions.

37.15. Informix Compatibility Mode

`ecpg` can be run in a so-called *Informix compatibility mode*. If this mode is active, it tries to behave as if it were the Informix precompiler for Informix E/SQL. Generally spoken this will allow you to use the dollar sign instead of the `EXEC SQL` primitive to introduce embedded SQL commands:

```
$int j = 3;
$CONNECT TO :dbname;
$CREATE TABLE test(i INT PRIMARY KEY, j INT);
$INSERT INTO test(i, j) VALUES (7, :j);
$COMMIT;
```

Note

There must not be any white space between the `$` and a following preprocessor directive, that is, `include`, `define`, `ifdef`, etc. Otherwise, the preprocessor will parse the token as a host variable.

There are two compatibility modes: `INFORMIX`, `INFORMIX_SE`

When linking programs that use this compatibility mode, remember to link against `libcompat` that is shipped with ECPG.

Besides the previously explained syntactic sugar, the Informix compatibility mode ports some functions for input, output and transformation of data as well as embedded SQL statements known from E/SQL to ECPG.

Informix compatibility mode is closely connected to the `pgtypeslib` library of ECPG. `pgtypeslib` maps SQL data types to data types within the C host program and most of the additional functions of the Informix compatibility mode allow you to operate on those C host program types. Note however that the extent of the compatibility is limited. It does not try to copy Informix behavior; it allows you to do more or less the same operations and gives you functions that have the same name and the same basic behavior but it is no drop-in replacement if you are using Informix at the moment. Moreover, some of the data types are different. For example, Postgres Pro's `datetime` and `interval` types do not know about ranges like for example `YEAR TO MINUTE` so you won't find support in ECPG for that either.

37.15.1. Additional Types

The Informix-special "string" pseudo-type for storing right-trimmed character string data is now supported in Informix-mode without using `typedef`. In fact, in Informix-mode, ECPG refuses to process source files that contain `typedef sometype string`;

```
EXEC SQL BEGIN DECLARE SECTION;
string userid; /* this variable will contain trimmed data */
EXEC SQL END DECLARE SECTION;
```

```
EXEC SQL FETCH MYCUR INTO :userid;
```

37.15.2. Additional/Missing Embedded SQL Statements

```
CLOSE DATABASE
```

This statement closes the current connection. In fact, this is a synonym for ECPG's `DISCONNECT CURRENT`:

```
$CLOSE DATABASE; /* close the current connection */
EXEC SQL CLOSE DATABASE;
```


FREE cursor_name

Due to differences in how ECPG works compared to Informix's ESQL/C (namely, which steps are purely grammar transformations and which steps rely on the underlying run-time library) there is no `FREE cursor_name` statement in ECPG. This is because in ECPG, `DECLARE CURSOR` doesn't translate to a function call into the run-time library that uses to the cursor name. This means that there's no run-time bookkeeping of SQL cursors in the ECPG run-time library, only in the Postgres Pro server.

FREE statement_name

`FREE statement_name` is a synonym for `DEALLOCATE PREPARE statement_name`.

37.15.3. Informix-compatible SQLDA Descriptor Areas

Informix-compatible mode supports a different structure than the one described in [Section 37.7.2](#). See below:

```
struct sqlvar_compat
{
    short    sqltype;
    int      sqllen;
    char     *sqldata;
    short    *sqlind;
    char     *sqlname;
    char     *sqlformat;
    short    sqlitype;
    short    sqlilen;
    char     *sqlidata;
    int      sqlxid;
    char     *sqltypename;
    short    sqltypelen;
    short    sqlownerlen;
    short    sqlsourcetype;
    char     *sqlownername;
    int      sqlsourceid;
    char     *sqlilongdata;
    int      sqlflags;
    void     *sqlreserved;
};

struct sqlda_compat
{
    short    sqld;
    struct sqlvar_compat *sqlvar;
    char     desc_name[19];
    short    desc_occ;
    struct sqlda_compat *desc_next;
    void     *reserved;
};

typedef struct sqlvar_compat    sqlvar_t;
typedef struct sqlda_compat    sqlda_t;
```

The global properties are:

`sqld`

The number of fields in the `SQLDA` descriptor.

`sqlvar`

Pointer to the per-field properties.

desc_name

Unused, filled with zero-bytes.

desc_occ

Size of the allocated structure.

desc_next

Pointer to the next SQLDA structure if the result set contains more than one record.

reserved

Unused pointer, contains NULL. Kept for Informix-compatibility.

The per-field properties are below, they are stored in the `sqlvar` array:

sqltype

Type of the field. Constants are in `sqltypes.h`

sqllen

Length of the field data.

sqldata

Pointer to the field data. The pointer is of `char *` type, the data pointed by it is in a binary format.
Example:

```
int intval;

switch (sqldata->sqlvar[i].sqltype)
{
    case SQLINTEGER:
        intval = *(int *)sqldata->sqlvar[i].sqldata;
        break;
    ...
}
```

sqlind

Pointer to the NULL indicator. If returned by DESCRIBE or FETCH then it's always a valid pointer. If used as input for EXECUTE ... USING sqlda; then NULL-pointer value means that the value for this field is non-NULL. Otherwise a valid pointer and `sqltype` has to be properly set. Example:

```
if (*(int2 *)sqldata->sqlvar[i].sqlind != 0)
    printf("value is NULL\n");
```

sqlname

Name of the field. 0-terminated string.

sqlformat

Reserved in Informix, value of `PQfformat` for the field.

sqlitype

Type of the NULL indicator data. It's always SQLSMINT when returning data from the server. When the SQLDA is used for a parameterized query, the data is treated according to the set type.

sqlilen

Length of the NULL indicator data.

sqlxid

Extended type of the field, result of [PQftype](#).

sqltypename
sqltypelen
sqlownerlen
sqlsourcetype
sqlownername
sqlsourceid
sqlflags
sqlreserved

Unused.

sqlilongdata

It equals to `sqldata` if `sqlllen` is larger than 32kB.

Example:

```
EXEC SQL INCLUDE sqlda.h;

sqlda_t          *sqlda; /* This doesn't need to be under embedded DECLARE SECTION */

EXEC SQL BEGIN DECLARE SECTION;
char *prep_stmt = "select * from table1";
int i;
EXEC SQL END DECLARE SECTION;

...

EXEC SQL PREPARE mystmt FROM :prep_stmt;

EXEC SQL DESCRIBE mystmt INTO sqlda;

printf("# of fields: %d\n", sqlda->sqld);
for (i = 0; i < sqlda->sqld; i++)
    printf("field %d: \"%s\"\n", sqlda->sqlvar[i]->sqlname);

EXEC SQL DECLARE mycursor CURSOR FOR mystmt;
EXEC SQL OPEN mycursor;
EXEC SQL WHENEVER NOT FOUND GOTO out;

while (1)
{
    EXEC SQL FETCH mycursor USING sqlda;
}

EXEC SQL CLOSE mycursor;

free(sqlda); /* The main structure is all to be free(),
              * sqlda and sqlda->sqlvar is in one allocated area */
```

37.15.4. Additional Functions

decadd

Add two decimal type values.

```
int decadd(decimal *arg1, decimal *arg2, decimal *sum);
```

The function receives a pointer to the first operand of type decimal (*arg1*), a pointer to the second operand of type decimal (*arg2*) and a pointer to a value of type decimal that will contain the sum (*sum*). On success, the function returns 0. `ECPG_INFORMIX_NUM_OVERFLOW` is returned in case of overflow and `ECPG_INFORMIX_NUM_UNDERFLOW` in case of underflow. -1 is returned for other failures and `errno` is set to the respective `errno` number of the `pgtypeslib`.

`deccmp`

Compare two variables of type decimal.

```
int deccmp(decimal *arg1, decimal *arg2);
```

The function receives a pointer to the first decimal value (*arg1*), a pointer to the second decimal value (*arg2*) and returns an integer value that indicates which is the bigger value.

- 1, if the value that *arg1* points to is bigger than the value that *arg2* points to
- -1, if the value that *arg1* points to is smaller than the value that *arg2* points to
- 0, if the value that *arg1* points to and the value that *arg2* points to are equal

`deccopy`

Copy a decimal value.

```
void deccopy(decimal *src, decimal *target);
```

The function receives a pointer to the decimal value that should be copied as the first argument (*src*) and a pointer to the target structure of type decimal (*target*) as the second argument.

`deccvasc`

Convert a value from its ASCII representation into a decimal type.

```
int deccvasc(char *cp, int len, decimal *np);
```

The function receives a pointer to string that contains the string representation of the number to be converted (*cp*) as well as its length *len*. *np* is a pointer to the decimal value that saves the result of the operation.

Valid formats are for example: -2, .794, +3.44, 592.49E07 or -32.84e-4.

The function returns 0 on success. If overflow or underflow occurred, `ECPG_INFORMIX_NUM_OVERFLOW` or `ECPG_INFORMIX_NUM_UNDERFLOW` is returned. If the ASCII representation could not be parsed, `ECPG_INFORMIX_BAD_NUMERIC` is returned or `ECPG_INFORMIX_BAD_EXPONENT` if this problem occurred while parsing the exponent.

`deccvdbl`

Convert a value of type double to a value of type decimal.

```
int deccvdbl(double dbl, decimal *np);
```

The function receives the variable of type double that should be converted as its first argument (*dbl*). As the second argument (*np*), the function receives a pointer to the decimal variable that should hold the result of the operation.

The function returns 0 on success and a negative value if the conversion failed.

`deccvint`

Convert a value of type int to a value of type decimal.

```
int deccvint(int in, decimal *np);
```

The function receives the variable of type int that should be converted as its first argument (*in*). As the second argument (*np*), the function receives a pointer to the decimal variable that should hold the result of the operation.

The function returns 0 on success and a negative value if the conversion failed.

deccvlong

Convert a value of type long to a value of type decimal.

```
int deccvlong(long lng, decimal *np);
```

The function receives the variable of type long that should be converted as its first argument (*lng*). As the second argument (*np*), the function receives a pointer to the decimal variable that should hold the result of the operation.

The function returns 0 on success and a negative value if the conversion failed.

decdiv

Divide two variables of type decimal.

```
int decdiv(decimal *n1, decimal *n2, decimal *result);
```

The function receives pointers to the variables that are the first (*n1*) and the second (*n2*) operands and calculates *n1/n2*. *result* is a pointer to the variable that should hold the result of the operation.

On success, 0 is returned and a negative value if the division fails. If overflow or underflow occurred, the function returns `ECPG_INFORMIX_NUM_OVERFLOW` or `ECPG_INFORMIX_NUM_UNDERFLOW` respectively. If an attempt to divide by zero is observed, the function returns `ECPG_INFORMIX_DIVIDE_ZERO`.

decmul

Multiply two decimal values.

```
int decmul(decimal *n1, decimal *n2, decimal *result);
```

The function receives pointers to the variables that are the first (*n1*) and the second (*n2*) operands and calculates *n1*n2*. *result* is a pointer to the variable that should hold the result of the operation.

On success, 0 is returned and a negative value if the multiplication fails. If overflow or underflow occurred, the function returns `ECPG_INFORMIX_NUM_OVERFLOW` or `ECPG_INFORMIX_NUM_UNDERFLOW` respectively.

decsub

Subtract one decimal value from another.

```
int decsub(decimal *n1, decimal *n2, decimal *result);
```

The function receives pointers to the variables that are the first (*n1*) and the second (*n2*) operands and calculates *n1-n2*. *result* is a pointer to the variable that should hold the result of the operation.

On success, 0 is returned and a negative value if the subtraction fails. If overflow or underflow occurred, the function returns `ECPG_INFORMIX_NUM_OVERFLOW` or `ECPG_INFORMIX_NUM_UNDERFLOW` respectively.

dectoasc

Convert a variable of type decimal to its ASCII representation in a C char* string.

```
int dectoasc(decimal *np, char *cp, int len, int right)
```

The function receives a pointer to a variable of type decimal (*np*) that it converts to its textual representation. *cp* is the buffer that should hold the result of the operation. The parameter *right* specifies, how many digits right of the decimal point should be included in the output. The result will be rounded to this number of decimal digits. Setting *right* to -1 indicates that all available decimal digits should be included in the output. If the length of the output buffer, which is indicated by *len* is not sufficient to hold the textual representation including the trailing zero byte, only a single * character is stored in the result and -1 is returned.

The function returns either -1 if the buffer `cp` was too small or `ECPG_INFORMIX_OUT_OF_MEMORY` if memory was exhausted.

dectodbl

Convert a variable of type decimal to a double.

```
int dectodbl(decimal *np, double *dblp);
```

The function receives a pointer to the decimal value to convert (`np`) and a pointer to the double variable that should hold the result of the operation (`dblp`).

On success, 0 is returned and a negative value if the conversion failed.

dectoint

Convert a variable of type decimal to an integer.

```
int dectoint(decimal *np, int *ip);
```

The function receives a pointer to the decimal value to convert (`np`) and a pointer to the integer variable that should hold the result of the operation (`ip`).

On success, 0 is returned and a negative value if the conversion failed. If an overflow occurred, `ECPG_INFORMIX_NUM_OVERFLOW` is returned.

Note that the ECPG implementation differs from the Informix implementation. Informix limits an integer to the range from -32767 to 32767, while the limits in the ECPG implementation depend on the architecture (`INT_MIN` .. `INT_MAX`).

dectolong

Convert a variable of type decimal to a long integer.

```
int dectolong(decimal *np, long *lngp);
```

The function receives a pointer to the decimal value to convert (`np`) and a pointer to the long variable that should hold the result of the operation (`lngp`).

On success, 0 is returned and a negative value if the conversion failed. If an overflow occurred, `ECPG_INFORMIX_NUM_OVERFLOW` is returned.

Note that the ECPG implementation differs from the Informix implementation. Informix limits a long integer to the range from -2,147,483,647 to 2,147,483,647, while the limits in the ECPG implementation depend on the architecture (`-LONG_MAX` .. `LONG_MAX`).

rdatestr

Converts a date to a C char* string.

```
int rdatestr(date d, char *str);
```

The function receives two arguments, the first one is the date to convert (`d`) and the second one is a pointer to the target string. The output format is always `yyyy-mm-dd`, so you need to allocate at least 11 bytes (including the zero-byte terminator) for the string.

The function returns 0 on success and a negative value in case of error.

Note that ECPG's implementation differs from the Informix implementation. In Informix the format can be influenced by setting environment variables. In ECPG however, you cannot change the output format.

rstrdate

Parse the textual representation of a date.

```
int rstrdate(char *str, date *d);
```

The function receives the textual representation of the date to convert (`str`) and a pointer to a variable of type `date` (`d`). This function does not allow you to specify a format mask. It uses the default format mask of Informix which is `mm/dd/yyyy`. Internally, this function is implemented by means of `rdefmtdate`. Therefore, `rstrdate` is not faster and if you have the choice you should opt for `rdefmtdate` which allows you to specify the format mask explicitly.

The function returns the same values as `rdefmtdate`.

`rtoday`

Get the current date.

```
void rtoday(date *d);
```

The function receives a pointer to a date variable (`d`) that it sets to the current date.

Internally this function uses the `PGTYPESdate_today` function.

`rjulmdy`

Extract the values for the day, the month and the year from a variable of type `date`.

```
int rjulmdy(date d, short mdy[3]);
```

The function receives the date `d` and a pointer to an array of 3 short integer values `mdy`. The variable name indicates the sequential order: `mdy[0]` will be set to contain the number of the month, `mdy[1]` will be set to the value of the day and `mdy[2]` will contain the year.

The function always returns 0 at the moment.

Internally the function uses the `PGTYPESdate_julmdy` function.

`rdefmtdate`

Use a format mask to convert a character string to a value of type `date`.

```
int rdefmtdate(date *d, char *fmt, char *str);
```

The function receives a pointer to the date value that should hold the result of the operation (`d`), the format mask to use for parsing the date (`fmt`) and the C `char*` string containing the textual representation of the date (`str`). The textual representation is expected to match the format mask. However you do not need to have a 1:1 mapping of the string to the format mask. The function only analyzes the sequential order and looks for the literals `yy` or `yyyy` that indicate the position of the year, `mm` to indicate the position of the month and `dd` to indicate the position of the day.

The function returns the following values:

- 0 - The function terminated successfully.
- `ECPG_INFORMIX_ENOSHORTDATE` - The date does not contain delimiters between day, month and year. In this case the input string must be exactly 6 or 8 bytes long but isn't.
- `ECPG_INFORMIX_ENOTDMY` - The format string did not correctly indicate the sequential order of year, month and day.
- `ECPG_INFORMIX_BAD_DAY` - The input string does not contain a valid day.
- `ECPG_INFORMIX_BAD_MONTH` - The input string does not contain a valid month.
- `ECPG_INFORMIX_BAD_YEAR` - The input string does not contain a valid year.

Internally this function is implemented to use the `PGTYPESdate_defmt_asc` function. See the reference there for a table of example input.

`rfmtdate`

Convert a variable of type `date` to its textual representation using a format mask.

```
int rfmtdate(date d, char *fmt, char *str);
```

The function receives the date to convert (*d*), the format mask (*fmt*) and the string that will hold the textual representation of the date (*str*).

On success, 0 is returned and a negative value if an error occurred.

Internally this function uses the [PGTYPESdate_fmt_asc](#) function, see the reference there for examples.

rmcyjul

Create a date value from an array of 3 short integers that specify the day, the month and the year of the date.

```
int rmcyjul(short mdy[3], date *d);
```

The function receives the array of the 3 short integers (*mdy*) and a pointer to a variable of type *date* that should hold the result of the operation.

Currently the function returns always 0.

Internally the function is implemented to use the function [PGTYPESdate_mdcyjul](#).

rdayofweek

Return a number representing the day of the week for a date value.

```
int rdayofweek(date d);
```

The function receives the date variable *d* as its only argument and returns an integer that indicates the day of the week for this date.

- 0 - Sunday
- 1 - Monday
- 2 - Tuesday
- 3 - Wednesday
- 4 - Thursday
- 5 - Friday
- 6 - Saturday

Internally the function is implemented to use the function [PGTYPESdate_dayofweek](#).

dtcurrent

Retrieve the current timestamp.

```
void dtcurrent(timestamp *ts);
```

The function retrieves the current timestamp and saves it into the timestamp variable that *ts* points to.

dtcvasc

Parses a timestamp from its textual representation into a timestamp variable.

```
int dtcvasc(char *str, timestamp *ts);
```

The function receives the string to parse (*str*) and a pointer to the timestamp variable that should hold the result of the operation (*ts*).

The function returns 0 on success and a negative value in case of error.

Internally this function uses the [PGTYPEStimestamp_from_asc](#) function. See the reference there for a table with example inputs.

dtcvfmtasc

Parses a timestamp from its textual representation using a format mask into a timestamp variable.

```
dtcvfmtasc(char *inbuf, char *fmtstr, timestamp *dtvalue)
```

The function receives the string to parse (*inbuf*), the format mask to use (*fmtstr*) and a pointer to the timestamp variable that should hold the result of the operation (*dtvalue*).

This function is implemented by means of the `PGTYPEStimestamp_defmt_asc` function. See the documentation there for a list of format specifiers that can be used.

The function returns 0 on success and a negative value in case of error.

dtsub

Subtract one timestamp from another and return a variable of type interval.

```
int dtsub(timestamp *ts1, timestamp *ts2, interval *iv);
```

The function will subtract the timestamp variable that *ts2* points to from the timestamp variable that *ts1* points to and will store the result in the interval variable that *iv* points to.

Upon success, the function returns 0 and a negative value if an error occurred.

dttoasc

Convert a timestamp variable to a C char* string.

```
int dttoasc(timestamp *ts, char *output);
```

The function receives a pointer to the timestamp variable to convert (*ts*) and the string that should hold the result of the operation (*output*). It converts *ts* to its textual representation according to the SQL standard, which is be YYYY-MM-DD HH:MM:SS.

Upon success, the function returns 0 and a negative value if an error occurred.

dttofmtasc

Convert a timestamp variable to a C char* using a format mask.

```
int dttofmtasc(timestamp *ts, char *output, int str_len, char *fmtstr);
```

The function receives a pointer to the timestamp to convert as its first argument (*ts*), a pointer to the output buffer (*output*), the maximal length that has been allocated for the output buffer (*str_len*) and the format mask to use for the conversion (*fmtstr*).

Upon success, the function returns 0 and a negative value if an error occurred.

Internally, this function uses the `PGTYPEStimestamp_fmt_asc` function. See the reference there for information on what format mask specifiers can be used.

intoasc

Convert an interval variable to a C char* string.

```
int intoasc(interval *i, char *str);
```

The function receives a pointer to the interval variable to convert (*i*) and the string that should hold the result of the operation (*str*). It converts *i* to its textual representation according to the SQL standard, which is be YYYY-MM-DD HH:MM:SS.

Upon success, the function returns 0 and a negative value if an error occurred.

rfmtlong

Convert a long integer value to its textual representation using a format mask.

```
int rfmtlong(long lng_val, char *fmt, char *outbuf);
```

The function receives the long value `lng_val`, the format mask `fmt` and a pointer to the output buffer `outbuf`. It converts the long value according to the format mask to its textual representation.

The format mask can be composed of the following format specifying characters:

- `*` (asterisk) - if this position would be blank otherwise, fill it with an asterisk.
- `&` (ampersand) - if this position would be blank otherwise, fill it with a zero.
- `#` - turn leading zeroes into blanks.
- `<` - left-justify the number in the string.
- `,` (comma) - group numbers of four or more digits into groups of three digits separated by a comma.
- `.` (period) - this character separates the whole-number part of the number from the fractional part.
- `-` (minus) - the minus sign appears if the number is a negative value.
- `+` (plus) - the plus sign appears if the number is a positive value.
- `(` - this replaces the minus sign in front of the negative number. The minus sign will not appear.
- `)` - this character replaces the minus and is printed behind the negative value.
- `$` - the currency symbol.

`rupshift`

Convert a string to upper case.

```
void rupshift(char *str);
```

The function receives a pointer to the string and transforms every lower case character to upper case.

`byleng`

Return the number of characters in a string without counting trailing blanks.

```
int byleng(char *str, int len);
```

The function expects a fixed-length string as its first argument (`str`) and its length as its second argument (`len`). It returns the number of significant characters, that is the length of the string without trailing blanks.

`ldchar`

Copy a fixed-length string into a null-terminated string.

```
void ldchar(char *src, int len, char *dest);
```

The function receives the fixed-length string to copy (`src`), its length (`len`) and a pointer to the destination memory (`dest`). Note that you need to reserve at least `len+1` bytes for the string that `dest` points to. The function copies at most `len` bytes to the new location (less if the source string has trailing blanks) and adds the null-terminator.

`rgetmsg`

```
int rgetmsg(int msgnum, char *s, int maxsize);
```

This function exists but is not implemented at the moment!

`rtypalign`

```
int rtypalign(int offset, int type);
```

This function exists but is not implemented at the moment!

rtypmsize

```
int rtypmsize(int type, int len);
```

This function exists but is not implemented at the moment!

rtypwidth

```
int rtypwidth(int sqltype, int sqllen);
```

This function exists but is not implemented at the moment!

rsetnull

Set a variable to NULL.

```
int rsetnull(int t, char *ptr);
```

The function receives an integer that indicates the type of the variable and a pointer to the variable itself that is cast to a C char* pointer.

The following types exist:

- CCHARTYPE - For a variable of type char or char*
- CSHORTTYPE - For a variable of type short int
- CINTTYPE - For a variable of type int
- CBOOLTYPE - For a variable of type boolean
- CFLOATTYPE - For a variable of type float
- CLONGTYPE - For a variable of type long
- CDOUBLETTYPE - For a variable of type double
- CDECIMALTYPE - For a variable of type decimal
- CDATETYPE - For a variable of type date
- CDTIMETYPE - For a variable of type timestamp

Here is an example of a call to this function:

```
$char c[] = "abc          ";
$short s = 17;
$int i = -74874;

rsetnull(CCHARTYPE, (char *) c);
rsetnull(CSHORTTYPE, (char *) &s);
rsetnull(CINTTYPE, (char *) &i);
```

risnull

Test if a variable is NULL.

```
int risnull(int t, char *ptr);
```

The function receives the type of the variable to test (`t`) as well a pointer to this variable (`ptr`). Note that the latter needs to be cast to a char*. See the function [rsetnull](#) for a list of possible variable types.

Here is an example of how to use this function:

```
$char c[] = "abc          ";
$short s = 17;
$int i = -74874;
```

```
risnull(CCHARTYPE, (char *) c);  
risnull(CSHORTTYPE, (char *) &s);  
risnull(CINTTYPE, (char *) &i);
```

37.15.5. Additional Constants

Note that all constants here describe errors and all of them are defined to represent negative values. In the descriptions of the different constants you can also find the value that the constants represent in the current implementation. However you should not rely on this number. You can however rely on the fact all of them are defined to represent negative values.

ECPG_INFORMIX_NUM_OVERFLOW

Functions return this value if an overflow occurred in a calculation. Internally it is defined as -1200 (the Informix definition).

ECPG_INFORMIX_NUM_UNDERFLOW

Functions return this value if an underflow occurred in a calculation. Internally it is defined as -1201 (the Informix definition).

ECPG_INFORMIX_DIVIDE_ZERO

Functions return this value if an attempt to divide by zero is observed. Internally it is defined as -1202 (the Informix definition).

ECPG_INFORMIX_BAD_YEAR

Functions return this value if a bad value for a year was found while parsing a date. Internally it is defined as -1204 (the Informix definition).

ECPG_INFORMIX_BAD_MONTH

Functions return this value if a bad value for a month was found while parsing a date. Internally it is defined as -1205 (the Informix definition).

ECPG_INFORMIX_BAD_DAY

Functions return this value if a bad value for a day was found while parsing a date. Internally it is defined as -1206 (the Informix definition).

ECPG_INFORMIX_ENOSHORTDATE

Functions return this value if a parsing routine needs a short date representation but did not get the date string in the right length. Internally it is defined as -1209 (the Informix definition).

ECPG_INFORMIX_DATE_CONVERT

Functions return this value if an error occurred during date formatting. Internally it is defined as -1210 (the Informix definition).

ECPG_INFORMIX_OUT_OF_MEMORY

Functions return this value if memory was exhausted during their operation. Internally it is defined as -1211 (the Informix definition).

ECPG_INFORMIX_ENOTDMY

Functions return this value if a parsing routine was supposed to get a format mask (like `mmddy`) but not all fields were listed correctly. Internally it is defined as -1212 (the Informix definition).

ECPG_INFORMIX_BAD_NUMERIC

Functions return this value either if a parsing routine cannot parse the textual representation for a numeric value because it contains errors or if a routine cannot complete a calculation involving

numeric variables because at least one of the numeric variables is invalid. Internally it is defined as -1213 (the Informix definition).

ECPG_INFORMIX_BAD_EXPONENT

Functions return this value if a parsing routine cannot parse an exponent. Internally it is defined as -1216 (the Informix definition).

ECPG_INFORMIX_BAD_DATE

Functions return this value if a parsing routine cannot parse a date. Internally it is defined as -1218 (the Informix definition).

ECPG_INFORMIX_EXTRA_CHARS

Functions return this value if a parsing routine is passed extra characters it cannot parse. Internally it is defined as -1264 (the Informix definition).

37.16. Oracle Compatibility Mode

`ecpg` can be run in a so-called *Oracle compatibility mode*. If this mode is active, it tries to behave as if it were Oracle Pro*C.

Specifically, this mode changes `ecpg` in three ways:

- Pad character arrays receiving character string types with trailing spaces to the specified length
- Zero byte terminate these character arrays, and set the indicator variable if truncation occurs
- Set the null indicator to -1 when character arrays receive empty character string types

37.17. Internals

This section explains how ECPG works internally. This information can occasionally be useful to help users understand how to use ECPG.

The first four lines written by `ecpg` to the output are fixed lines. Two are comments and two are include lines necessary to interface to the library. Then the preprocessor reads through the file and writes output. Normally it just echoes everything to the output.

When it sees an `EXEC SQL` statement, it intervenes and changes it. The command starts with `EXEC SQL` and ends with `;`. Everything in between is treated as an SQL statement and parsed for variable substitution.

Variable substitution occurs when a symbol starts with a colon (`:`). The variable with that name is looked up among the variables that were previously declared within a `EXEC SQL DECLARE` section.

The most important function in the library is `ECPGdo`, which takes care of executing most commands. It takes a variable number of arguments. This can easily add up to 50 or so arguments, and we hope this will not be a problem on any platform.

The arguments are:

A line number

This is the line number of the original line; used in error messages only.

A string

This is the SQL command that is to be issued. It is modified by the input variables, i.e., the variables that were not known at compile time but are to be entered in the command. Where the variables should go the string contains `?`.

Input variables

Every input variable causes ten arguments to be created. (See below.)

ECPGt_EOIT

An enum telling that there are no more input variables.

Output variables

Every output variable causes ten arguments to be created. (See below.) These variables are filled by the function.

ECPGt_EORT

An enum telling that there are no more variables.

For every variable that is part of the SQL command, the function gets ten arguments:

1. The type as a special symbol.
2. A pointer to the value or a pointer to the pointer.
3. The size of the variable if it is a `char` or `varchar`.
4. The number of elements in the array (for array fetches).
5. The offset to the next element in the array (for array fetches).
6. The type of the indicator variable as a special symbol.
7. A pointer to the indicator variable.
8. 0
9. The number of elements in the indicator array (for array fetches).
10. The offset to the next element in the indicator array (for array fetches).

Note that not all SQL commands are treated in this way. For instance, an open cursor statement like:

```
EXEC SQL OPEN cursor;
```

is not copied to the output. Instead, the cursor's `DECLARE` command is used at the position of the `OPEN` command because it indeed opens the cursor.

Here is a complete example describing the output of the preprocessor of a file `foo.pgc` (details might change with each particular version of the preprocessor):

```
EXEC SQL BEGIN DECLARE SECTION;
int index;
int result;
EXEC SQL END DECLARE SECTION;
...
EXEC SQL SELECT res INTO :result FROM mytable WHERE index = :index;
```

is translated into:

```
/* Processed by ecpg (2.6.0) */
/* These two include files are added by the preprocessor */
#include <ecpgtype.h>;
#include <ecpglib.h>;

/* exec sql begin declare section */

#line 1 "foo.pgc"

    int index;
    int result;
/* exec sql end declare section */
...
ECPGdo(__LINE__, NULL, "SELECT res FROM mytable WHERE index = ?      ",
```

```
    ECPGt_int,&(index),1L,1L,sizeof(int),
    ECPGt_NO_INDICATOR, NULL , 0L, 0L, 0L, ECPGt_EOIT,
    ECPGt_int,&(result),1L,1L,sizeof(int),
    ECPGt_NO_INDICATOR, NULL , 0L, 0L, 0L, ECPGt_EORT);
#line 147 "foo.pgc"
```

(The indentation here is added for readability and not something the preprocessor does.)

Chapter 38. The Information Schema

The information schema consists of a set of views that contain information about the objects defined in the current database. The information schema is defined in the SQL standard and can therefore be expected to be portable and remain stable — unlike the system catalogs, which are specific to Postgres Pro and are modeled after implementation concerns. The information schema views do not, however, contain information about Postgres Pro-specific features; to inquire about those you need to query the system catalogs or other Postgres Pro-specific views.

Note

When querying the database for constraint information, it is possible for a standard-compliant query that expects to return one row to return several. This is because the SQL standard requires constraint names to be unique within a schema, but Postgres Pro does not enforce this restriction. Postgres Pro automatically-generated constraint names avoid duplicates in the same schema, but users can specify such duplicate names.

This problem can appear when querying information schema views such as `check_constraint_routine_usage`, `check_constraints`, `domain_constraints`, and `referential_constraints`. Some other views have similar issues but contain the table name to help distinguish duplicate rows, e.g., `constraint_column_usage`, `constraint_table_usage`, `table_constraints`.

38.1. The Schema

The information schema itself is a schema named `information_schema`. This schema automatically exists in all databases. The owner of this schema is the initial database user in the cluster, and that user naturally has all the privileges on this schema, including the ability to drop it (but the space savings achieved by that are minuscule).

By default, the information schema is not in the schema search path, so you need to access all objects in it through qualified names. Since the names of some of the objects in the information schema are generic names that might occur in user applications, you should be careful if you want to put the information schema in the path.

38.2. Data Types

The columns of the information schema views use special data types that are defined in the information schema. These are defined as simple domains over ordinary built-in types. You should not use these types for work outside the information schema, but your applications must be prepared for them if they select from the information schema.

These types are:

`cardinal_number`

A nonnegative integer.

`character_data`

A character string (without specific maximum length).

`sql_identifier`

A character string. This type is used for SQL identifiers, the type `character_data` is used for any other kind of text data.

`time_stamp`

A domain over the type `timestamp with time zone`

`yes_or_no`

A character string domain that contains either `YES` or `NO`. This is used to represent Boolean (`true/false`) data in the information schema. (The information schema was invented before the type `boolean` was added to the SQL standard, so this convention is necessary to keep the information schema backward compatible.)

Every column in the information schema has one of these five types.

38.3. `information_schema_catalog_name`

`information_schema_catalog_name` is a table that always contains one row and one column containing the name of the current database (current catalog, in SQL terminology).

Table 38.1. `information_schema_catalog_name` Columns

Column Type	Description
<code>catalog_name sql_identifier</code>	Name of the database that contains this information schema

38.4. `administrable_role_authorizations`

The view `administrable_role_authorizations` identifies all roles that the current user has the admin option for.

Table 38.2. `administrable_role_authorizations` Columns

Column Type	Description
<code>grantee sql_identifier</code>	Name of the role to which this role membership was granted (can be the current user, or a different role in case of nested role memberships)
<code>role_name sql_identifier</code>	Name of a role
<code>is_grantable yes_or_no</code>	Always <code>YES</code>

38.5. `applicable_roles`

The view `applicable_roles` identifies all roles whose privileges the current user can use. This means there is some chain of role grants from the current user to the role in question. The current user itself is also an applicable role. The set of applicable roles is generally used for permission checking.

Table 38.3. `applicable_roles` Columns

Column Type	Description
<code>grantee sql_identifier</code>	Name of the role to which this role membership was granted (can be the current user, or a different role in case of nested role memberships)
<code>role_name sql_identifier</code>	Name of a role
<code>is_grantable yes_or_no</code>	<code>YES</code> if the grantee has the admin option on the role, <code>NO</code> if not

38.6. attributes

The view `attributes` contains information about the attributes of composite data types defined in the database. (Note that the view does not give information about table columns, which are sometimes called attributes in Postgres Pro contexts.) Only those attributes are shown that the current user has access to (by way of being the owner of or having some privilege on the type).

Table 38.4. attributes Columns

Column Type	Description
<code>udt_catalog sql_identifier</code>	Name of the database containing the data type (always the current database)
<code>udt_schema sql_identifier</code>	Name of the schema containing the data type
<code>udt_name sql_identifier</code>	Name of the data type
<code>attribute_name sql_identifier</code>	Name of the attribute
<code>ordinal_position cardinal_number</code>	Ordinal position of the attribute within the data type (count starts at 1)
<code>attribute_default character_data</code>	Default expression of the attribute
<code>is_nullable yes_or_no</code>	YES if the attribute is possibly nullable, NO if it is known not nullable.
<code>data_type character_data</code>	Data type of the attribute, if it is a built-in type, or ARRAY if it is some array (in that case, see the view <code>element_types</code>), else USER-DEFINED (in that case, the type is identified in <code>attribute_udt_name</code> and associated columns).
<code>character_maximum_length cardinal_number</code>	If <code>data_type</code> identifies a character or bit string type, the declared maximum length; null for all other data types or if no maximum length was declared.
<code>character_octet_length cardinal_number</code>	If <code>data_type</code> identifies a character type, the maximum possible length in octets (bytes) of a datum; null for all other data types. The maximum octet length depends on the declared character maximum length (see above) and the server encoding.
<code>character_set_catalog sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog sql_identifier</code>	Name of the database containing the collation of the attribute (always the current database), null if default or the data type of the attribute is not collatable
<code>collation_schema sql_identifier</code>	Name of the schema containing the collation of the attribute, null if default or the data type of the attribute is not collatable
<code>collation_name sql_identifier</code>	Name of the collation of the attribute, null if default or the data type of the attribute is not collatable

Column Type	Description
<code>numeric_precision</code> <code>cardinal_number</code>	If <code>data_type</code> identifies a numeric type, this column contains the (declared or implicit) precision of the type for this attribute. The precision indicates the number of significant digits. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>numeric_precision_radix</code> <code>cardinal_number</code>	If <code>data_type</code> identifies a numeric type, this column indicates in which base the values in the columns <code>numeric_precision</code> and <code>numeric_scale</code> are expressed. The value is either 2 or 10. For all other data types, this column is null.
<code>numeric_scale</code> <code>cardinal_number</code>	If <code>data_type</code> identifies an exact numeric type, this column contains the (declared or implicit) scale of the type for this attribute. The scale indicates the number of significant digits to the right of the decimal point. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>datetime_precision</code> <code>cardinal_number</code>	If <code>data_type</code> identifies a date, time, timestamp, or interval type, this column contains the (declared or implicit) fractional seconds precision of the type for this attribute, that is, the number of decimal digits maintained following the decimal point in the seconds value. For all other data types, this column is null.
<code>interval_type</code> <code>character_data</code>	If <code>data_type</code> identifies an interval type, this column contains the specification which fields the intervals include for this attribute, e.g., YEAR TO MONTH, DAY TO SECOND, etc. If no field restrictions were specified (that is, the interval accepts all fields), and for all other data types, this field is null.
<code>interval_precision</code> <code>cardinal_number</code>	Applies to a feature not available in Postgres Pro (see <code>datetime_precision</code> for the fractional seconds precision of interval type attributes)
<code>attribute_udt_catalog</code> <code>sql_identifier</code>	Name of the database that the attribute data type is defined in (always the current database)
<code>attribute_udt_schema</code> <code>sql_identifier</code>	Name of the schema that the attribute data type is defined in
<code>attribute_udt_name</code> <code>sql_identifier</code>	Name of the attribute data type
<code>scope_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_schema</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>maximum_cardinality</code> <code>cardinal_number</code>	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
<code>dtd_identifier</code> <code>sql_identifier</code>	An identifier of the data type descriptor of the attribute, unique among the data type descriptors pertaining to the composite type. This is mainly useful for joining with other instances of such identifiers. (The specific format of the identifier is not defined and not guaranteed to remain the same in future versions.)
<code>is_derived_reference_attribute</code> <code>yes_or_no</code>	Applies to a feature not available in Postgres Pro

See also under [Section 38.17](#), a similarly structured view, for further information on some of the columns.

38.7. character_sets

The view `character_sets` identifies the character sets available in the current database. Since Postgres Pro does not support multiple character sets within one database, this view only shows one, which is the database encoding.

Take note of how the following terms are used in the SQL standard:

character repertoire

An abstract collection of characters, for example `UNICODE`, `UCS`, or `LATIN1`. Not exposed as an SQL object, but visible in this view.

character encoding form

An encoding of some character repertoire. Most older character repertoires only use one encoding form, and so there are no separate names for them (e.g., `LATIN2` is an encoding form applicable to the `LATIN2` repertoire). But for example Unicode has the encoding forms `UTF8`, `UTF16`, etc. (not all supported by Postgres Pro Shardman). Encoding forms are not exposed as an SQL object, but are visible in this view.

character set

A named SQL object that identifies a character repertoire, a character encoding, and a default collation. A predefined character set would typically have the same name as an encoding form, but users could define other names. For example, the character set `UTF8` would typically identify the character repertoire `UCS`, encoding form `UTF8`, and some default collation.

You can think of an “encoding” in Postgres Pro either as a character set or a character encoding form. They will have the same name, and there can only be one in one database.

Table 38.5. character_sets Columns

Column Type	Description
<code>character_set_catalog</code> <code>sql_identifier</code>	Character sets are currently not implemented as schema objects, so this column is null.
<code>character_set_schema</code> <code>sql_identifier</code>	Character sets are currently not implemented as schema objects, so this column is null.
<code>character_set_name</code> <code>sql_identifier</code>	Name of the character set, currently implemented as showing the name of the database encoding
<code>character_repertoire</code> <code>sql_identifier</code>	Character repertoire, showing <code>UCS</code> if the encoding is <code>UTF8</code> , else just the encoding name
<code>form_of_use</code> <code>sql_identifier</code>	Character encoding form, same as the database encoding
<code>default_collate_catalog</code> <code>sql_identifier</code>	Name of the database containing the default collation (always the current database, if any collation is identified)
<code>default_collate_schema</code> <code>sql_identifier</code>	Name of the schema containing the default collation
<code>default_collate_name</code> <code>sql_identifier</code>	Name of the default collation. The default collation is identified as the collation that matches the <code>COLLATE</code> and <code>CTYPE</code> settings of the current database. If there is no such collation, then this column and the associated schema and catalog columns are null.

38.8. `check_constraint_routine_usage`

The view `check_constraint_routine_usage` identifies routines (functions and procedures) that are used by a check constraint. Only those routines are shown that are owned by a currently enabled role.

Table 38.6. `check_constraint_routine_usage` Columns

Column Type	Description
<code>constraint_catalog sql_identifier</code>	Name of the database containing the constraint (always the current database)
<code>constraint_schema sql_identifier</code>	Name of the schema containing the constraint
<code>constraint_name sql_identifier</code>	Name of the constraint
<code>specific_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.

38.9. `check_constraints`

The view `check_constraints` contains all check constraints, either defined on a table or on a domain, that are owned by a currently enabled role. (The owner of the table or domain is the owner of the constraint.)

The SQL standard considers not-null constraints to be check constraints with a `CHECK (column_name IS NOT NULL)` expression. So not-null constraints are also included here and don't have a separate view.

Table 38.7. `check_constraints` Columns

Column Type	Description
<code>constraint_catalog sql_identifier</code>	Name of the database containing the constraint (always the current database)
<code>constraint_schema sql_identifier</code>	Name of the schema containing the constraint
<code>constraint_name sql_identifier</code>	Name of the constraint
<code>check_clause character_data</code>	The check expression of the check constraint

38.10. `collations`

The view `collations` contains the collations available in the current database.

Table 38.8. `collations` Columns

Column Type	Description
<code>collation_catalog sql_identifier</code>	

Column Type	Description
	Name of the database containing the collation (always the current database)
collation_schema sql_identifier	Name of the schema containing the collation
collation_name sql_identifier	Name of the default collation
pad_attribute character_data	Always NO PAD (The alternative PAD SPACE is not supported by Postgres Pro.)

38.11. collation_character_set_applicability

The view `collation_character_set_applicability` identifies which character set the available collations are applicable to. In Postgres Pro, there is only one character set per database (see explanation in [Section 38.7](#)), so this view does not provide much useful information.

Table 38.9. collation_character_set_applicability Columns

Column Type	Description
	Name of the database containing the collation (always the current database)
collation_schema sql_identifier	Name of the schema containing the collation
collation_name sql_identifier	Name of the default collation
character_set_catalog sql_identifier	Character sets are currently not implemented as schema objects, so this column is null
character_set_schema sql_identifier	Character sets are currently not implemented as schema objects, so this column is null
character_set_name sql_identifier	Name of the character set

38.12. column_column_usage

The view `column_column_usage` identifies all generated columns that depend on another base column in the same table. Only tables owned by a currently enabled role are included.

Table 38.10. column_column_usage Columns

Column Type	Description
	Name of the database containing the table (always the current database)
table_schema sql_identifier	Name of the schema containing the table
table_name sql_identifier	Name of the table
column_name sql_identifier	Name of the base column that a generated column depends on
dependent_column sql_identifier	Name of the generated column

38.13. column_domain_usage

The view `column_domain_usage` identifies all columns (of a table or a view) that make use of some domain defined in the current database and owned by a currently enabled role.

Table 38.11. column_domain_usage Columns

Column Type	Description
domain_catalog sql_identifier	Name of the database containing the domain (always the current database)
domain_schema sql_identifier	Name of the schema containing the domain
domain_name sql_identifier	Name of the domain
table_catalog sql_identifier	Name of the database containing the table (always the current database)
table_schema sql_identifier	Name of the schema containing the table
table_name sql_identifier	Name of the table
column_name sql_identifier	Name of the column

38.14. column_options

The view `column_options` contains all the options defined for foreign table columns in the current database. Only those foreign table columns are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.12. column_options Columns

Column Type	Description
table_catalog sql_identifier	Name of the database that contains the foreign table (always the current database)
table_schema sql_identifier	Name of the schema that contains the foreign table
table_name sql_identifier	Name of the foreign table
column_name sql_identifier	Name of the column
option_name sql_identifier	Name of an option
option_value character_data	Value of the option

38.15. column_privileges

The view `column_privileges` identifies all privileges granted on columns to a currently enabled role or by a currently enabled role. There is one row for each combination of column, grantor, and grantee.

If a privilege has been granted on an entire table, it will show up in this view as a grant for each column, but only for the privilege types where column granularity is possible: `SELECT`, `INSERT`, `UPDATE`, `REFERENCES`.

Table 38.13. `column_privileges` Columns

Column Type	Description
<code>grantor sql_identifier</code>	Name of the role that granted the privilege
<code>grantee sql_identifier</code>	Name of the role that the privilege was granted to
<code>table_catalog sql_identifier</code>	Name of the database that contains the table that contains the column (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema that contains the table that contains the column
<code>table_name sql_identifier</code>	Name of the table that contains the column
<code>column_name sql_identifier</code>	Name of the column
<code>privilege_type character_data</code>	Type of the privilege: <code>SELECT</code> , <code>INSERT</code> , <code>UPDATE</code> , or <code>REFERENCES</code>
<code>is_grantable yes_or_no</code>	YES if the privilege is grantable, NO if not

38.16. `column_udt_usage`

The view `column_udt_usage` identifies all columns that use data types owned by a currently enabled role. Note that in Postgres Pro, built-in data types behave like user-defined types, so they are included here as well. See also [Section 38.17](#) for details.

Table 38.14. `column_udt_usage` Columns

Column Type	Description
<code>udt_catalog sql_identifier</code>	Name of the database that the column data type (the underlying type of the domain, if applicable) is defined in (always the current database)
<code>udt_schema sql_identifier</code>	Name of the schema that the column data type (the underlying type of the domain, if applicable) is defined in
<code>udt_name sql_identifier</code>	Name of the column data type (the underlying type of the domain, if applicable)
<code>table_catalog sql_identifier</code>	Name of the database containing the table (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema containing the table
<code>table_name sql_identifier</code>	Name of the table
<code>column_name sql_identifier</code>	Name of the column

38.17. columns

The view `columns` contains information about all table columns (or view columns) in the database. System columns (`ctid`, etc.) are not included. Only those columns are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.15. columns Columns

Column Type	Description
<code>table_catalog sql_identifier</code>	Name of the database containing the table (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema containing the table
<code>table_name sql_identifier</code>	Name of the table
<code>column_name sql_identifier</code>	Name of the column
<code>ordinal_position cardinal_number</code>	Ordinal position of the column within the table (count starts at 1)
<code>column_default character_data</code>	Default expression of the column
<code>is_nullable yes_or_no</code>	YES if the column is possibly nullable, NO if it is known not nullable. A not-null constraint is one way a column can be known not nullable, but there can be others.
<code>data_type character_data</code>	Data type of the column, if it is a built-in type, or ARRAY if it is some array (in that case, see the view <code>element_types</code>), else USER-DEFINED (in that case, the type is identified in <code>udt_name</code> and associated columns). If the column is based on a domain, this column refers to the type underlying the domain (and the domain is identified in <code>domain_name</code> and associated columns).
<code>character_maximum_length cardinal_number</code>	If <code>data_type</code> identifies a character or bit string type, the declared maximum length; null for all other data types or if no maximum length was declared.
<code>character_octet_length cardinal_number</code>	If <code>data_type</code> identifies a character type, the maximum possible length in octets (bytes) of a datum; null for all other data types. The maximum octet length depends on the declared character maximum length (see above) and the server encoding.
<code>numeric_precision cardinal_number</code>	If <code>data_type</code> identifies a numeric type, this column contains the (declared or implicit) precision of the type for this column. The precision indicates the number of significant digits. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>numeric_precision_radix cardinal_number</code>	If <code>data_type</code> identifies a numeric type, this column indicates in which base the values in the columns <code>numeric_precision</code> and <code>numeric_scale</code> are expressed. The value is either 2 or 10. For all other data types, this column is null.
<code>numeric_scale cardinal_number</code>	If <code>data_type</code> identifies an exact numeric type, this column contains the (declared or implicit) scale of the type for this column. The scale indicates the number of significant digits to the right of the decimal point. It can be expressed in decimal (base 10) or binary (base 2) terms,

Column Type	Description
	as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>datetime_precision</code> <code>cardinal_number</code>	If <code>data_type</code> identifies a date, time, timestamp, or interval type, this column contains the (declared or implicit) fractional seconds precision of the type for this column, that is, the number of decimal digits maintained following the decimal point in the seconds value. For all other data types, this column is null.
<code>interval_type</code> <code>character_data</code>	If <code>data_type</code> identifies an interval type, this column contains the specification which fields the intervals include for this column, e.g., YEAR TO MONTH, DAY TO SECOND, etc. If no field restrictions were specified (that is, the interval accepts all fields), and for all other data types, this field is null.
<code>interval_precision</code> <code>cardinal_number</code>	Applies to a feature not available in Postgres Pro (see <code>datetime_precision</code> for the fractional seconds precision of interval type columns)
<code>character_set_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog</code> <code>sql_identifier</code>	Name of the database containing the collation of the column (always the current database), null if default or the data type of the column is not collatable
<code>collation_schema</code> <code>sql_identifier</code>	Name of the schema containing the collation of the column, null if default or the data type of the column is not collatable
<code>collation_name</code> <code>sql_identifier</code>	Name of the collation of the column, null if default or the data type of the column is not collatable
<code>domain_catalog</code> <code>sql_identifier</code>	If the column has a domain type, the name of the database that the domain is defined in (always the current database), else null.
<code>domain_schema</code> <code>sql_identifier</code>	If the column has a domain type, the name of the schema that the domain is defined in, else null.
<code>domain_name</code> <code>sql_identifier</code>	If the column has a domain type, the name of the domain, else null.
<code>udt_catalog</code> <code>sql_identifier</code>	Name of the database that the column data type (the underlying type of the domain, if applicable) is defined in (always the current database)
<code>udt_schema</code> <code>sql_identifier</code>	Name of the schema that the column data type (the underlying type of the domain, if applicable) is defined in
<code>udt_name</code> <code>sql_identifier</code>	Name of the column data type (the underlying type of the domain, if applicable)
<code>scope_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro

Column Type	Description
<code>scope_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>maximum_cardinality cardinal_number</code>	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
<code>dtd_identifier sql_identifier</code>	An identifier of the data type descriptor of the column, unique among the data type descriptors pertaining to the table. This is mainly useful for joining with other instances of such identifiers. (The specific format of the identifier is not defined and not guaranteed to remain the same in future versions.)
<code>is_self_referencing yes_or_no</code>	Applies to a feature not available in Postgres Pro
<code>is_identity yes_or_no</code>	If the column is an identity column, then YES, else NO.
<code>identity_generation character_data</code>	If the column is an identity column, then ALWAYS or BY DEFAULT, reflecting the definition of the column.
<code>identity_start character_data</code>	If the column is an identity column, then the start value of the internal sequence, else null.
<code>identity_increment character_data</code>	If the column is an identity column, then the increment of the internal sequence, else null.
<code>identity_maximum character_data</code>	If the column is an identity column, then the maximum value of the internal sequence, else null.
<code>identity_minimum character_data</code>	If the column is an identity column, then the minimum value of the internal sequence, else null.
<code>identity_cycle yes_or_no</code>	If the column is an identity column, then YES if the internal sequence cycles or NO if it does not; otherwise null.
<code>is_generated character_data</code>	If the column is a generated column, then ALWAYS, else NEVER.
<code>generation_expression character_data</code>	If the column is a generated column, then the generation expression, else null.
<code>is_updatable yes_or_no</code>	YES if the column is updatable, NO if not (Columns in base tables are always updatable, columns in views not necessarily)

Since data types can be defined in a variety of ways in SQL, and Postgres Pro contains additional ways to define data types, their representation in the information schema can be somewhat difficult. The column `data_type` is supposed to identify the underlying built-in type of the column. In Postgres Pro, this means that the type is defined in the system catalog schema `pg_catalog`. This column might be useful if the application can handle the well-known built-in types specially (for example, format the numeric types differently or use the data in the precision columns). The columns `udt_name`, `udt_schema`, and `udt_catalog` always identify the underlying data type of the column, even if the column is based on a domain. (Since Postgres Pro treats built-in types like user-defined types, built-in types appear here as well. This is an extension of the SQL standard.) These columns should be used if an application wants to process data differently according to the type, because in that case it wouldn't matter if the column is

really based on a domain. If the column is based on a domain, the identity of the domain is stored in the columns `domain_name`, `domain_schema`, and `domain_catalog`. If you want to pair up columns with their associated data types and treat domains as separate types, you could write `coalesce(domain_name, udt_name)`, etc.

38.18. `constraint_column_usage`

The view `constraint_column_usage` identifies all columns in the current database that are used by some constraint. Only those columns are shown that are contained in a table owned by a currently enabled role. For a check constraint, this view identifies the columns that are used in the check expression. For a not-null constraint, this view identifies the column that the constraint is defined on. For a foreign key constraint, this view identifies the columns that the foreign key references. For a unique or primary key constraint, this view identifies the constrained columns.

Table 38.16. `constraint_column_usage` Columns

Column Type	Description
<code>table_catalog sql_identifier</code>	Name of the database that contains the table that contains the column that is used by some constraint (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema that contains the table that contains the column that is used by some constraint
<code>table_name sql_identifier</code>	Name of the table that contains the column that is used by some constraint
<code>column_name sql_identifier</code>	Name of the column that is used by some constraint
<code>constraint_catalog sql_identifier</code>	Name of the database that contains the constraint (always the current database)
<code>constraint_schema sql_identifier</code>	Name of the schema that contains the constraint
<code>constraint_name sql_identifier</code>	Name of the constraint

38.19. `constraint_table_usage`

The view `constraint_table_usage` identifies all tables in the current database that are used by some constraint and are owned by a currently enabled role. (This is different from the view `table_constraints`, which identifies all table constraints along with the table they are defined on.) For a foreign key constraint, this view identifies the table that the foreign key references. For a unique or primary key constraint, this view simply identifies the table the constraint belongs to. Check constraints and not-null constraints are not included in this view.

Table 38.17. `constraint_table_usage` Columns

Column Type	Description
<code>table_catalog sql_identifier</code>	Name of the database that contains the table that is used by some constraint (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema that contains the table that is used by some constraint
<code>table_name sql_identifier</code>	Name of the table that is used by some constraint

Column Type	Description
<code>constraint_catalog sql_identifier</code>	Name of the database that contains the constraint (always the current database)
<code>constraint_schema sql_identifier</code>	Name of the schema that contains the constraint
<code>constraint_name sql_identifier</code>	Name of the constraint

38.20. data_type_privileges

The view `data_type_privileges` identifies all data type descriptors that the current user has access to, by way of being the owner of the described object or having some privilege for it. A data type descriptor is generated whenever a data type is used in the definition of a table column, a domain, or a function (as parameter or return type) and stores some information about how the data type is used in that instance (for example, the declared maximum length, if applicable). Each data type descriptor is assigned an arbitrary identifier that is unique among the data type descriptor identifiers assigned for one object (table, domain, function). This view is probably not useful for applications, but it is used to define some other views in the information schema.

Table 38.18. data_type_privileges Columns

Column Type	Description
<code>object_catalog sql_identifier</code>	Name of the database that contains the described object (always the current database)
<code>object_schema sql_identifier</code>	Name of the schema that contains the described object
<code>object_name sql_identifier</code>	Name of the described object
<code>object_type character_data</code>	The type of the described object: one of <code>TABLE</code> (the data type descriptor pertains to a column of that table), <code>DOMAIN</code> (the data type descriptors pertains to that domain), <code>ROUTINE</code> (the data type descriptor pertains to a parameter or the return data type of that function).
<code>dtd_identifier sql_identifier</code>	The identifier of the data type descriptor, which is unique among the data type descriptors for that same object.

38.21. domain_constraints

The view `domain_constraints` contains all constraints belonging to domains defined in the current database. Only those domains are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.19. domain_constraints Columns

Column Type	Description
<code>constraint_catalog sql_identifier</code>	Name of the database that contains the constraint (always the current database)
<code>constraint_schema sql_identifier</code>	Name of the schema that contains the constraint
<code>constraint_name sql_identifier</code>	Name of the constraint

Column Type	Description
domain_catalog sql_identifier	Name of the database that contains the domain (always the current database)
domain_schema sql_identifier	Name of the schema that contains the domain
domain_name sql_identifier	Name of the domain
is_deferrable yes_or_no	YES if the constraint is deferrable, NO if not
initially_deferred yes_or_no	YES if the constraint is deferrable and initially deferred, NO if not

38.22. domain_udt_usage

The view `domain_udt_usage` identifies all domains that are based on data types owned by a currently enabled role. Note that in Postgres Pro, built-in data types behave like user-defined types, so they are included here as well.

Table 38.20. domain_udt_usage Columns

Column Type	Description
udt_catalog sql_identifier	Name of the database that the domain data type is defined in (always the current database)
udt_schema sql_identifier	Name of the schema that the domain data type is defined in
udt_name sql_identifier	Name of the domain data type
domain_catalog sql_identifier	Name of the database that contains the domain (always the current database)
domain_schema sql_identifier	Name of the schema that contains the domain
domain_name sql_identifier	Name of the domain

38.23. domains

The view `domains` contains all *domains* defined in the current database. Only those domains are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.21. domains Columns

Column Type	Description
domain_catalog sql_identifier	Name of the database that contains the domain (always the current database)
domain_schema sql_identifier	Name of the schema that contains the domain
domain_name sql_identifier	Name of the domain
data_type character_data	

Column Type	Description
	Data type of the domain, if it is a built-in type, or ARRAY if it is some array (in that case, see the view <code>element_types</code>), else USER-DEFINED (in that case, the type is identified in <code>udt_name</code> and associated columns).
<code>character_maximum_length</code> <code>cardinal_number</code>	If the domain has a character or bit string type, the declared maximum length; null for all other data types or if no maximum length was declared.
<code>character_octet_length</code> <code>cardinal_number</code>	If the domain has a character type, the maximum possible length in octets (bytes) of a datum; null for all other data types. The maximum octet length depends on the declared character maximum length (see above) and the server encoding.
<code>character_set_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog</code> <code>sql_identifier</code>	Name of the database containing the collation of the domain (always the current database), null if default or the data type of the domain is not collatable
<code>collation_schema</code> <code>sql_identifier</code>	Name of the schema containing the collation of the domain, null if default or the data type of the domain is not collatable
<code>collation_name</code> <code>sql_identifier</code>	Name of the collation of the domain, null if default or the data type of the domain is not collatable
<code>numeric_precision</code> <code>cardinal_number</code>	If the domain has a numeric type, this column contains the (declared or implicit) precision of the type for this domain. The precision indicates the number of significant digits. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>numeric_precision_radix</code> <code>cardinal_number</code>	If the domain has a numeric type, this column indicates in which base the values in the columns <code>numeric_precision</code> and <code>numeric_scale</code> are expressed. The value is either 2 or 10. For all other data types, this column is null.
<code>numeric_scale</code> <code>cardinal_number</code>	If the domain has an exact numeric type, this column contains the (declared or implicit) scale of the type for this domain. The scale indicates the number of significant digits to the right of the decimal point. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> . For all other data types, this column is null.
<code>datetime_precision</code> <code>cardinal_number</code>	If <code>data_type</code> identifies a date, time, timestamp, or interval type, this column contains the (declared or implicit) fractional seconds precision of the type for this domain, that is, the number of decimal digits maintained following the decimal point in the seconds value. For all other data types, this column is null.
<code>interval_type</code> <code>character_data</code>	If <code>data_type</code> identifies an interval type, this column contains the specification which fields the intervals include for this domain, e.g., YEAR TO MONTH, DAY TO SECOND, etc. If no field restrictions were specified (that is, the interval accepts all fields), and for all other data types, this field is null.

Column Type	Description
<code>interval_precision</code> <code>cardinal_number</code>	Applies to a feature not available in Postgres Pro (see <code>datetime_precision</code> for the fractional seconds precision of interval type domains)
<code>domain_default</code> <code>character_data</code>	Default expression of the domain
<code>udt_catalog</code> <code>sql_identifier</code>	Name of the database that the domain data type is defined in (always the current database)
<code>udt_schema</code> <code>sql_identifier</code>	Name of the schema that the domain data type is defined in
<code>udt_name</code> <code>sql_identifier</code>	Name of the domain data type
<code>scope_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_schema</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>maximum_cardinality</code> <code>cardinal_number</code>	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
<code>dtd_identifier</code> <code>sql_identifier</code>	An identifier of the data type descriptor of the domain, unique among the data type descriptors pertaining to the domain (which is trivial, because a domain only contains one data type descriptor). This is mainly useful for joining with other instances of such identifiers. (The specific format of the identifier is not defined and not guaranteed to remain the same in future versions.)

38.24. `element_types`

The view `element_types` contains the data type descriptors of the elements of arrays. When a table column, composite-type attribute, domain, function parameter, or function return value is defined to be of an array type, the respective information schema view only contains `ARRAY` in the column `data_type`. To obtain information on the element type of the array, you can join the respective view with this view. For example, to show the columns of a table with data types and array element types, if applicable, you could do:

```
SELECT c.column_name, c.data_type, e.data_type AS element_type
FROM information_schema.columns c LEFT JOIN information_schema.element_types e
    ON ((c.table_catalog, c.table_schema, c.table_name, 'TABLE', c.dtd_identifier)
        = (e.object_catalog, e.object_schema, e.object_name, e.object_type,
            e.collection_type_identifier))
WHERE c.table_schema = '...' AND c.table_name = '...'
ORDER BY c.ordinal_position;
```

This view only includes objects that the current user has access to, by way of being the owner or having some privilege.

Table 38.22. `element_types` Columns

Column Type	Description
<code>object_catalog</code> <code>sql_identifier</code>	Name of the database that contains the object that uses the array being described (always the current database)

Column Type	Description
object_schema sql_identifier	Name of the schema that contains the object that uses the array being described
object_name sql_identifier	Name of the object that uses the array being described
object_type character_data	The type of the object that uses the array being described: one of TABLE (the array is used by a column of that table), USER-DEFINED TYPE (the array is used by an attribute of that composite type), DOMAIN (the array is used by that domain), ROUTINE (the array is used by a parameter or the return data type of that function).
collection_type_identifier sql_identifier	The identifier of the data type descriptor of the array being described. Use this to join with the dtd_identifier columns of other information schema views.
data_type character_data	Data type of the array elements, if it is a built-in type, else USER-DEFINED (in that case, the type is identified in udt_name and associated columns).
character_maximum_length cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
character_octet_length cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
character_set_catalog sql_identifier	Applies to a feature not available in Postgres Pro
character_set_schema sql_identifier	Applies to a feature not available in Postgres Pro
character_set_name sql_identifier	Applies to a feature not available in Postgres Pro
collation_catalog sql_identifier	Name of the database containing the collation of the element type (always the current database), null if default or the data type of the element is not collatable
collation_schema sql_identifier	Name of the schema containing the collation of the element type, null if default or the data type of the element is not collatable
collation_name sql_identifier	Name of the collation of the element type, null if default or the data type of the element is not collatable
numeric_precision cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
numeric_precision_radix cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
numeric_scale cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
datetime_precision cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
interval_type character_data	Always null, since this information is not applied to array element data types in Postgres Pro
interval_precision cardinal_number	Always null, since this information is not applied to array element data types in Postgres Pro
udt_catalog sql_identifier	

Column Type	Description
	Name of the database that the data type of the elements is defined in (always the current database)
udt_schema sql_identifier	Name of the schema that the data type of the elements is defined in
udt_name sql_identifier	Name of the data type of the elements
scope_catalog sql_identifier	Applies to a feature not available in Postgres Pro
scope_schema sql_identifier	Applies to a feature not available in Postgres Pro
scope_name sql_identifier	Applies to a feature not available in Postgres Pro
maximum_cardinality cardinal_number	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
dtd_identifier sql_identifier	An identifier of the data type descriptor of the element. This is currently not useful.

38.25. enabled_roles

The view `enabled_roles` identifies the currently “enabled roles”. The enabled roles are recursively defined as the current user together with all roles that have been granted to the enabled roles with automatic inheritance. In other words, these are all roles that the current user has direct or indirect, automatically inheriting membership in.

For permission checking, the set of “applicable roles” is applied, which can be broader than the set of enabled roles. So generally, it is better to use the view `applicable_roles` instead of this one; See [Section 38.5](#) for details on `applicable_roles` view.

Table 38.23. enabled_roles Columns

Column Type	Description
role_name sql_identifier	Name of a role

38.26. foreign_data_wrapper_options

The view `foreign_data_wrapper_options` contains all the options defined for foreign-data wrappers in the current database. Only those foreign-data wrappers are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.24. foreign_data_wrapper_options Columns

Column Type	Description
foreign_data_wrapper_catalog sql_identifier	Name of the database that the foreign-data wrapper is defined in (always the current database)
foreign_data_wrapper_name sql_identifier	Name of the foreign-data wrapper
option_name sql_identifier	Name of an option

Column Type	Description
option_value character_data	Value of the option

38.27. foreign_data_wrappers

The view `foreign_data_wrappers` contains all foreign-data wrappers defined in the current database. Only those foreign-data wrappers are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.25. foreign_data_wrappers Columns

Column Type	Description
foreign_data_wrapper_catalog sql_identifier	Name of the database that contains the foreign-data wrapper (always the current database)
foreign_data_wrapper_name sql_identifier	Name of the foreign-data wrapper
authorization_identifier sql_identifier	Name of the owner of the foreign server
library_name character_data	File name of the library that implementing this foreign-data wrapper
foreign_data_wrapper_language character_data	Language used to implement this foreign-data wrapper

38.28. foreign_server_options

The view `foreign_server_options` contains all the options defined for foreign servers in the current database. Only those foreign servers are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.26. foreign_server_options Columns

Column Type	Description
foreign_server_catalog sql_identifier	Name of the database that the foreign server is defined in (always the current database)
foreign_server_name sql_identifier	Name of the foreign server
option_name sql_identifier	Name of an option
option_value character_data	Value of the option

38.29. foreign_servers

The view `foreign_servers` contains all foreign servers defined in the current database. Only those foreign servers are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.27. foreign_servers Columns

Column Type	Description
foreign_server_catalog sql_identifier	

Column Type	Description
	Name of the database that the foreign server is defined in (always the current database)
foreign_server_name sql_identifier	Name of the foreign server
foreign_data_wrapper_catalog sql_identifier	Name of the database that contains the foreign-data wrapper used by the foreign server (always the current database)
foreign_data_wrapper_name sql_identifier	Name of the foreign-data wrapper used by the foreign server
foreign_server_type character_data	Foreign server type information, if specified upon creation
foreign_server_version character_data	Foreign server version information, if specified upon creation
authorization_identifier sql_identifier	Name of the owner of the foreign server

38.30. foreign_table_options

The view `foreign_table_options` contains all the options defined for foreign tables in the current database. Only those foreign tables are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.28. foreign_table_options Columns

Column Type	Description
foreign_table_catalog sql_identifier	Name of the database that contains the foreign table (always the current database)
foreign_table_schema sql_identifier	Name of the schema that contains the foreign table
foreign_table_name sql_identifier	Name of the foreign table
option_name sql_identifier	Name of an option
option_value character_data	Value of the option

38.31. foreign_tables

The view `foreign_tables` contains all foreign tables defined in the current database. Only those foreign tables are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.29. foreign_tables Columns

Column Type	Description
foreign_table_catalog sql_identifier	Name of the database that the foreign table is defined in (always the current database)
foreign_table_schema sql_identifier	Name of the schema that contains the foreign table
foreign_table_name sql_identifier	Name of the foreign table

Column Type	Description
foreign_server_catalog sql_identifier	Name of the database that the foreign server is defined in (always the current database)
foreign_server_name sql_identifier	Name of the foreign server

38.32. key_column_usage

The view `key_column_usage` identifies all columns in the current database that are restricted by some unique, primary key, or foreign key constraint. Check constraints are not included in this view. Only those columns are shown that the current user has access to, by way of being the owner or having some privilege.

Table 38.30. key_column_usage Columns

Column Type	Description
constraint_catalog sql_identifier	Name of the database that contains the constraint (always the current database)
constraint_schema sql_identifier	Name of the schema that contains the constraint
constraint_name sql_identifier	Name of the constraint
table_catalog sql_identifier	Name of the database that contains the table that contains the column that is restricted by this constraint (always the current database)
table_schema sql_identifier	Name of the schema that contains the table that contains the column that is restricted by this constraint
table_name sql_identifier	Name of the table that contains the column that is restricted by this constraint
column_name sql_identifier	Name of the column that is restricted by this constraint
ordinal_position cardinal_number	Ordinal position of the column within the constraint key (count starts at 1)
position_in_unique_constraint cardinal_number	For a foreign-key constraint, ordinal position of the referenced column within its unique constraint (count starts at 1); otherwise null

38.33. parameters

The view `parameters` contains information about the parameters (arguments) of all functions in the current database. Only those functions are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.31. parameters Columns

Column Type	Description
specific_catalog sql_identifier	Name of the database containing the function (always the current database)
specific_schema sql_identifier	Name of the schema containing the function

Column Type	Description
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.
<code>ordinal_position cardinal_number</code>	Ordinal position of the parameter in the argument list of the function (count starts at 1)
<code>parameter_mode character_data</code>	IN for input parameter, OUT for output parameter, and INOUT for input/output parameter.
<code>is_result yes_or_no</code>	Applies to a feature not available in Postgres Pro
<code>as_locator yes_or_no</code>	Applies to a feature not available in Postgres Pro
<code>parameter_name sql_identifier</code>	Name of the parameter, or null if the parameter has no name
<code>data_type character_data</code>	Data type of the parameter, if it is a built-in type, or ARRAY if it is some array (in that case, see the view <code>element_types</code>), else USER-DEFINED (in that case, the type is identified in <code>udt_name</code> and associated columns).
<code>character_maximum_length cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>character_octet_length cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>character_set_catalog sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog sql_identifier</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>collation_schema sql_identifier</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>collation_name sql_identifier</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>numeric_precision cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>numeric_precision_radix cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>numeric_scale cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>datetime_precision cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>interval_type character_data</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>interval_precision cardinal_number</code>	Always null, since this information is not applied to parameter data types in Postgres Pro
<code>udt_catalog sql_identifier</code>	

Column Type	Description
	Name of the database that the data type of the parameter is defined in (always the current database)
udt_schema sql_identifier	Name of the schema that the data type of the parameter is defined in
udt_name sql_identifier	Name of the data type of the parameter
scope_catalog sql_identifier	Applies to a feature not available in Postgres Pro
scope_schema sql_identifier	Applies to a feature not available in Postgres Pro
scope_name sql_identifier	Applies to a feature not available in Postgres Pro
maximum_cardinality cardinal_number	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
dtd_identifier sql_identifier	An identifier of the data type descriptor of the parameter, unique among the data type descriptors pertaining to the function. This is mainly useful for joining with other instances of such identifiers. (The specific format of the identifier is not defined and not guaranteed to remain the same in future versions.)
parameter_default character_data	The default expression of the parameter, or null if none or if the function is not owned by a currently enabled role.

38.34. referential_constraints

The view `referential_constraints` contains all referential (foreign key) constraints in the current database. Only those constraints are shown for which the current user has write access to the referencing table (by way of being the owner or having some privilege other than `SELECT`).

Table 38.32. referential_constraints Columns

Column Type	Description
constraint_catalog sql_identifier	Name of the database containing the constraint (always the current database)
constraint_schema sql_identifier	Name of the schema containing the constraint
constraint_name sql_identifier	Name of the constraint
unique_constraint_catalog sql_identifier	Name of the database that contains the unique or primary key constraint that the foreign key constraint references (always the current database)
unique_constraint_schema sql_identifier	Name of the schema that contains the unique or primary key constraint that the foreign key constraint references
unique_constraint_name sql_identifier	Name of the unique or primary key constraint that the foreign key constraint references
match_option character_data	Match option of the foreign key constraint: <code>FULL</code> , <code>PARTIAL</code> , or <code>NONE</code> .
update_rule character_data	

Column Type	Description
	Update rule of the foreign key constraint: CASCADE, SET NULL, SET DEFAULT, RESTRICT, or NO ACTION.
delete_rule character_data	Delete rule of the foreign key constraint: CASCADE, SET NULL, SET DEFAULT, RESTRICT, or NO ACTION.

38.35. role_column_grants

The view `role_column_grants` identifies all privileges granted on columns where the grantor or grantee is a currently enabled role. Further information can be found under `column_privileges`. The only effective difference between this view and `column_privileges` is that this view omits columns that have been made accessible to the current user by way of a grant to `PUBLIC`.

Table 38.33. role_column_grants Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to
table_catalog sql_identifier	Name of the database that contains the table that contains the column (always the current database)
table_schema sql_identifier	Name of the schema that contains the table that contains the column
table_name sql_identifier	Name of the table that contains the column
column_name sql_identifier	Name of the column
privilege_type character_data	Type of the privilege: SELECT, INSERT, UPDATE, or REFERENCES
is_grantable yes_or_no	YES if the privilege is grantable, NO if not

38.36. role_routine_grants

The view `role_routine_grants` identifies all privileges granted on functions where the grantor or grantee is a currently enabled role. Further information can be found under `routine_privileges`. The only effective difference between this view and `routine_privileges` is that this view omits functions that have been made accessible to the current user by way of a grant to `PUBLIC`.

Table 38.34. role_routine_grants Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to
specific_catalog sql_identifier	

Column Type	Description
	Name of the database containing the function (always the current database)
specific_schema sql_identifier	Name of the schema containing the function
specific_name sql_identifier	The “specific name” of the function. See Section 38.45 for more information.
routine_catalog sql_identifier	Name of the database containing the function (always the current database)
routine_schema sql_identifier	Name of the schema containing the function
routine_name sql_identifier	Name of the function (might be duplicated in case of overloading)
privilege_type character_data	Always EXECUTE (the only privilege type for functions)
is_grantable yes_or_no	YES if the privilege is grantable, NO if not

38.37. role_table_grants

The view `role_table_grants` identifies all privileges granted on tables or views where the grantor or grantee is a currently enabled role. Further information can be found under `table_privileges`. The only effective difference between this view and `table_privileges` is that this view omits tables that have been made accessible to the current user by way of a grant to `PUBLIC`.

Table 38.35. role_table_grants Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to
table_catalog sql_identifier	Name of the database that contains the table (always the current database)
table_schema sql_identifier	Name of the schema that contains the table
table_name sql_identifier	Name of the table
privilege_type character_data	Type of the privilege: SELECT, INSERT, UPDATE, DELETE, TRUNCATE, REFERENCES, or TRIGGER
is_grantable yes_or_no	YES if the privilege is grantable, NO if not
with_hierarchy yes_or_no	In the SQL standard, WITH HIERARCHY OPTION is a separate (sub-)privilege allowing certain operations on table inheritance hierarchies. In Postgres Pro, this is included in the SELECT privilege, so this column shows YES if the privilege is SELECT, else NO.

38.38. role_udt_grants

The view `role_udt_grants` is intended to identify `USAGE` privileges granted on user-defined types where the grantor or grantee is a currently enabled role. Further information can be found under `udt_priv-`

ileges. The only effective difference between this view and `udt_privileges` is that this view omits objects that have been made accessible to the current user by way of a grant to `PUBLIC`. Since data types do not have real privileges in Postgres Pro, but only an implicit grant to `PUBLIC`, this view is empty.

Table 38.36. `role_udt_grants` Columns

Column Type	Description
<code>grantor sql_identifier</code>	The name of the role that granted the privilege
<code>grantee sql_identifier</code>	The name of the role that the privilege was granted to
<code>udt_catalog sql_identifier</code>	Name of the database containing the type (always the current database)
<code>udt_schema sql_identifier</code>	Name of the schema containing the type
<code>udt_name sql_identifier</code>	Name of the type
<code>privilege_type character_data</code>	Always <code>TYPE USAGE</code>
<code>is_grantable yes_or_no</code>	YES if the privilege is grantable, NO if not

38.39. `role_usage_grants`

The view `role_usage_grants` identifies `USAGE` privileges granted on various kinds of objects where the grantor or grantee is a currently enabled role. Further information can be found under `usage_privileges`. The only effective difference between this view and `usage_privileges` is that this view omits objects that have been made accessible to the current user by way of a grant to `PUBLIC`.

Table 38.37. `role_usage_grants` Columns

Column Type	Description
<code>grantor sql_identifier</code>	The name of the role that granted the privilege
<code>grantee sql_identifier</code>	The name of the role that the privilege was granted to
<code>object_catalog sql_identifier</code>	Name of the database containing the object (always the current database)
<code>object_schema sql_identifier</code>	Name of the schema containing the object, if applicable, else an empty string
<code>object_name sql_identifier</code>	Name of the object
<code>object_type character_data</code>	<code>COLLATION</code> or <code>DOMAIN</code> or <code>FOREIGN DATA WRAPPER</code> or <code>FOREIGN SERVER</code> or <code>SEQUENCE</code>
<code>privilege_type character_data</code>	Always <code>USAGE</code>
<code>is_grantable yes_or_no</code>	YES if the privilege is grantable, NO if not

38.40. routine_column_usage

The view `routine_column_usage` identifies all columns that are used by a function or procedure, either in the SQL body or in parameter default expressions. (This only works for unquoted SQL bodies, not quoted bodies or functions in other languages.) A column is only included if its table is owned by a currently enabled role.

Table 38.38. routine_column_usage Columns

Column Type	Description
<code>specific_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.
<code>routine_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>routine_schema sql_identifier</code>	Name of the schema containing the function
<code>routine_name sql_identifier</code>	Name of the function (might be duplicated in case of overloading)
<code>table_catalog sql_identifier</code>	Name of the database that contains the table that is used by the function (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema that contains the table that is used by the function
<code>table_name sql_identifier</code>	Name of the table that is used by the function
<code>column_name sql_identifier</code>	Name of the column that is used by the function

38.41. routine_privileges

The view `routine_privileges` identifies all privileges granted on functions to a currently enabled role or by a currently enabled role. There is one row for each combination of function, grantor, and grantee.

Table 38.39. routine_privileges Columns

Column Type	Description
<code>grantor sql_identifier</code>	Name of the role that granted the privilege
<code>grantee sql_identifier</code>	Name of the role that the privilege was granted to
<code>specific_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.
<code>routine_catalog sql_identifier</code>	

Column Type	Description
	Name of the database containing the function (always the current database)
<code>routine_schema</code> <code>sql_identifier</code>	Name of the schema containing the function
<code>routine_name</code> <code>sql_identifier</code>	Name of the function (might be duplicated in case of overloading)
<code>privilege_type</code> <code>character_data</code>	Always EXECUTE (the only privilege type for functions)
<code>is_grantable</code> <code>yes_or_no</code>	YES if the privilege is grantable, NO if not

38.42. `routine_routine_usage`

The view `routine_routine_usage` identifies all functions or procedures that are used by another (or the same) function or procedure, either in the SQL body or in parameter default expressions. (This only works for unquoted SQL bodies, not quoted bodies or functions in other languages.) An entry is included here only if the used function is owned by a currently enabled role. (There is no such restriction on the using function.)

Note that the entries for both functions in the view refer to the “specific” name of the routine, even though the column names are used in a way that is inconsistent with other information schema views about routines. This is per SQL standard, although it is arguably a misdesign. See [Section 38.45](#) for more information about specific names.

Table 38.40. `routine_routine_usage` Columns

Column Type	Description
<code>specific_catalog</code> <code>sql_identifier</code>	Name of the database containing the using function (always the current database)
<code>specific_schema</code> <code>sql_identifier</code>	Name of the schema containing the using function
<code>specific_name</code> <code>sql_identifier</code>	The “specific name” of the using function.
<code>routine_catalog</code> <code>sql_identifier</code>	Name of the database that contains the function that is used by the first function (always the current database)
<code>routine_schema</code> <code>sql_identifier</code>	Name of the schema that contains the function that is used by the first function
<code>routine_name</code> <code>sql_identifier</code>	The “specific name” of the function that is used by the first function.

38.43. `routine_sequence_usage`

The view `routine_sequence_usage` identifies all sequences that are used by a function or procedure, either in the SQL body or in parameter default expressions. (This only works for unquoted SQL bodies, not quoted bodies or functions in other languages.) A sequence is only included if that sequence is owned by a currently enabled role.

Table 38.41. `routine_sequence_usage` Columns

Column Type	Description
<code>specific_catalog</code> <code>sql_identifier</code>	

Column Type	Description
	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.
<code>routine_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>routine_schema sql_identifier</code>	Name of the schema containing the function
<code>routine_name sql_identifier</code>	Name of the function (might be duplicated in case of overloading)
<code>schema_catalog sql_identifier</code>	Name of the database that contains the sequence that is used by the function (always the current database)
<code>sequence_schema sql_identifier</code>	Name of the schema that contains the sequence that is used by the function
<code>sequence_name sql_identifier</code>	Name of the sequence that is used by the function

38.44. routine_table_usage

The view `routine_table_usage` is meant to identify all tables that are used by a function or procedure. This information is currently not tracked by Postgres Pro.

Table 38.42. `routine_table_usage` Columns

Column Type	Description
<code>specific_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. See Section 38.45 for more information.
<code>routine_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>routine_schema sql_identifier</code>	Name of the schema containing the function
<code>routine_name sql_identifier</code>	Name of the function (might be duplicated in case of overloading)
<code>table_catalog sql_identifier</code>	Name of the database that contains the table that is used by the function (always the current database)
<code>table_schema sql_identifier</code>	Name of the schema that contains the table that is used by the function
<code>table_name sql_identifier</code>	Name of the table that is used by the function

38.45. routines

The view `routines` contains all functions and procedures in the current database. Only those functions and procedures are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.43. routines Columns

Column Type	Description
<code>specific_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>specific_schema sql_identifier</code>	Name of the schema containing the function
<code>specific_name sql_identifier</code>	The “specific name” of the function. This is a name that uniquely identifies the function in the schema, even if the real name of the function is overloaded. The format of the specific name is not defined, it should only be used to compare it to other instances of specific routine names.
<code>routine_catalog sql_identifier</code>	Name of the database containing the function (always the current database)
<code>routine_schema sql_identifier</code>	Name of the schema containing the function
<code>routine_name sql_identifier</code>	Name of the function (might be duplicated in case of overloading)
<code>routine_type character_data</code>	FUNCTION for a function, PROCEDURE for a procedure
<code>module_catalog sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>module_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>module_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>udt_catalog sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>udt_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>udt_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>data_type character_data</code>	Return data type of the function, if it is a built-in type, or ARRAY if it is some array (in that case, see the view <code>element_types</code>), else USER-DEFINED (in that case, the type is identified in <code>type_udt_name</code> and associated columns). Null for a procedure.
<code>character_maximum_length cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>character_octet_length cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>character_set_catalog sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema sql_identifier</code>	Applies to a feature not available in Postgres Pro

Column Type	Description
<code>character_set_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog</code> <code>sql_identifier</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>collation_schema</code> <code>sql_identifier</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>collation_name</code> <code>sql_identifier</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>numeric_precision</code> <code>cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>numeric_precision_radix</code> <code>cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>numeric_scale</code> <code>cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>datetime_precision</code> <code>cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>interval_type</code> <code>character_data</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>interval_precision</code> <code>cardinal_number</code>	Always null, since this information is not applied to return data types in Postgres Pro
<code>type_udt_catalog</code> <code>sql_identifier</code>	Name of the database that the return data type of the function is defined in (always the current database). Null for a procedure.
<code>type_udt_schema</code> <code>sql_identifier</code>	Name of the schema that the return data type of the function is defined in. Null for a procedure.
<code>type_udt_name</code> <code>sql_identifier</code>	Name of the return data type of the function. Null for a procedure.
<code>scope_catalog</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_schema</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>scope_name</code> <code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>maximum_cardinality</code> <code>cardinal_number</code>	Always null, because arrays always have unlimited maximum cardinality in Postgres Pro
<code>dtd_identifier</code> <code>sql_identifier</code>	An identifier of the data type descriptor of the return data type of this function, unique among the data type descriptors pertaining to the function. This is mainly useful for joining with other instances of such identifiers. (The specific format of the identifier is not defined and not guaranteed to remain the same in future versions.)
<code>routine_body</code> <code>character_data</code>	If the function is an SQL function, then SQL, else EXTERNAL.
<code>routine_definition</code> <code>character_data</code>	The source text of the function (null if the function is not owned by a currently enabled role). (According to the SQL standard, this column is only applicable if <code>routine_body</code> is SQL, but in

Column Type	Description
	Postgres Pro it will contain whatever source text was specified when the function was created.)
external_name character_data	If this function is a C function, then the external name (link symbol) of the function; else null. (This works out to be the same value that is shown in routine_definition .)
external_language character_data	The language the function is written in
parameter_style character_data	Always GENERAL (The SQL standard defines other parameter styles, which are not available in Postgres Pro.)
is_deterministic yes_or_no	If the function is declared immutable (called deterministic in the SQL standard), then YES, else NO. (You cannot query the other volatility levels available in Postgres Pro through the information schema.)
sql_data_access character_data	Always MODIFIES, meaning that the function possibly modifies SQL data. This information is not useful for Postgres Pro.
is_null_call yes_or_no	If the function automatically returns null if any of its arguments are null, then YES, else NO. Null for a procedure.
sql_path character_data	Applies to a feature not available in Postgres Pro
schema_level_routine yes_or_no	Always YES (The opposite would be a method of a user-defined type, which is a feature not available in Postgres Pro.)
max_dynamic_result_sets cardinal_number	Applies to a feature not available in Postgres Pro
is_user_defined_cast yes_or_no	Applies to a feature not available in Postgres Pro
is_implicitly_invocable yes_or_no	Applies to a feature not available in Postgres Pro
security_type character_data	If the function runs with the privileges of the current user, then INVOKER, if the function runs with the privileges of the user who defined it, then DEFINER.
to_sql_specific_catalog sql_identifier	Applies to a feature not available in Postgres Pro
to_sql_specific_schema sql_identifier	Applies to a feature not available in Postgres Pro
to_sql_specific_name sql_identifier	Applies to a feature not available in Postgres Pro
as_locator yes_or_no	Applies to a feature not available in Postgres Pro
created_time_stamp	Applies to a feature not available in Postgres Pro
last_altered time_stamp	Applies to a feature not available in Postgres Pro
new_savepoint_level yes_or_no	

Column Type	Description
	Applies to a feature not available in Postgres Pro
is_udt_dependent yes_or_no	Currently always NO. The alternative YES applies to a feature not available in Postgres Pro.
result_cast_from_data_type character_data	Applies to a feature not available in Postgres Pro
result_cast_as_locator yes_or_no	Applies to a feature not available in Postgres Pro
result_cast_char_max_length cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_char_octet_length cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_char_set_catalog sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_char_set_schema sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_char_set_name sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_collation_catalog sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_collation_schema sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_collation_name sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_numeric_precision cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_numeric_precision_radix cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_numeric_scale cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_datetime_precision cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_interval_type character_data	Applies to a feature not available in Postgres Pro
result_cast_interval_precision cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_type_udt_catalog sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_type_udt_schema sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_type_udt_name sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_scope_catalog sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_scope_schema sql_identifier	Applies to a feature not available in Postgres Pro
result_cast_scope_name sql_identifier	

Column Type	Description
	Applies to a feature not available in Postgres Pro
result_cast_maximum_cardinality cardinal_number	Applies to a feature not available in Postgres Pro
result_cast_dtd_identifier sql_identifier	Applies to a feature not available in Postgres Pro

38.46. schemata

The view `schemata` contains all schemas in the current database that the current user has access to (by way of being the owner or having some privilege).

Table 38.44. schemata Columns

Column Type	Description
catalog_name sql_identifier	Name of the database that the schema is contained in (always the current database)
schema_name sql_identifier	Name of the schema
schema_owner sql_identifier	Name of the owner of the schema
default_character_set_catalog sql_identifier	Applies to a feature not available in Postgres Pro
default_character_set_schema sql_identifier	Applies to a feature not available in Postgres Pro
default_character_set_name sql_identifier	Applies to a feature not available in Postgres Pro
sql_path character_data	Applies to a feature not available in Postgres Pro

38.47. sequences

The view `sequences` contains all sequences defined in the current database. Only those sequences are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.45. sequences Columns

Column Type	Description
sequence_catalog sql_identifier	Name of the database that contains the sequence (always the current database)
sequence_schema sql_identifier	Name of the schema that contains the sequence
sequence_name sql_identifier	Name of the sequence
data_type character_data	The data type of the sequence.
numeric_precision cardinal_number	This column contains the (declared or implicit) precision of the sequence data type (see above). The precision indicates the number of significant digits. It can be expressed in dec-

Column Type	Description
	imal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> .
<code>numeric_precision_radix</code> <code>cardinal_number</code>	This column indicates in which base the values in the columns <code>numeric_precision</code> and <code>numeric_scale</code> are expressed. The value is either 2 or 10.
<code>numeric_scale</code> <code>cardinal_number</code>	This column contains the (declared or implicit) scale of the sequence data type (see above). The scale indicates the number of significant digits to the right of the decimal point. It can be expressed in decimal (base 10) or binary (base 2) terms, as specified in the column <code>numeric_precision_radix</code> .
<code>start_value</code> <code>character_data</code>	The start value of the sequence
<code>minimum_value</code> <code>character_data</code>	The minimum value of the sequence
<code>maximum_value</code> <code>character_data</code>	The maximum value of the sequence
<code>increment</code> <code>character_data</code>	The increment of the sequence
<code>cycle_option</code> <code>yes_or_no</code>	YES if the sequence cycles, else NO

Note that in accordance with the SQL standard, the start, minimum, maximum, and increment values are returned as character strings.

38.48. `sql_features`

The table `sql_features` contains information about which formal features defined in the SQL standard are supported by Postgres Pro. This is the same information that is presented in [Appendix D](#). There you can also find some additional background information.

Table 38.46. `sql_features` Columns

Column Type	Description
<code>feature_id</code> <code>character_data</code>	Identifier string of the feature
<code>feature_name</code> <code>character_data</code>	Descriptive name of the feature
<code>sub_feature_id</code> <code>character_data</code>	Identifier string of the subfeature, or a zero-length string if not a subfeature
<code>sub_feature_name</code> <code>character_data</code>	Descriptive name of the subfeature, or a zero-length string if not a subfeature
<code>is_supported</code> <code>yes_or_no</code>	YES if the feature is fully supported by the current version of Postgres Pro, NO if not
<code>is_verified_by</code> <code>character_data</code>	Always null, since the Postgres Pro development group does not perform formal testing of feature conformance
<code>comments</code> <code>character_data</code>	Possibly a comment about the supported status of the feature

38.49. sql_implementation_info

The table `sql_implementation_info` contains information about various aspects that are left implementation-defined by the SQL standard. This information is primarily intended for use in the context of the ODBC interface; users of other interfaces will probably find this information to be of little use. For this reason, the individual implementation information items are not described here; you will find them in the description of the ODBC interface.

Table 38.47. sql_implementation_info Columns

Column Type	Description
<code>implementation_info_id</code> <code>character_data</code>	Identifier string of the implementation information item
<code>implementation_info_name</code> <code>character_data</code>	Descriptive name of the implementation information item
<code>integer_value</code> <code>cardinal_number</code>	Value of the implementation information item, or null if the value is contained in the column <code>character_value</code>
<code>character_value</code> <code>character_data</code>	Value of the implementation information item, or null if the value is contained in the column <code>integer_value</code>
<code>comments</code> <code>character_data</code>	Possibly a comment pertaining to the implementation information item

38.50. sql_parts

The table `sql_parts` contains information about which of the several parts of the SQL standard are supported by Postgres Pro.

Table 38.48. sql_parts Columns

Column Type	Description
<code>feature_id</code> <code>character_data</code>	An identifier string containing the number of the part
<code>feature_name</code> <code>character_data</code>	Descriptive name of the part
<code>is_supported</code> <code>yes_or_no</code>	YES if the part is fully supported by the current version of Postgres Pro, NO if not
<code>is_verified_by</code> <code>character_data</code>	Always null, since the Postgres Pro development group does not perform formal testing of feature conformance
<code>comments</code> <code>character_data</code>	Possibly a comment about the supported status of the part

38.51. sql_sizing

The table `sql_sizing` contains information about various size limits and maximum values in Postgres Pro. This information is primarily intended for use in the context of the ODBC interface; users of other interfaces will probably find this information to be of little use. For this reason, the individual sizing items are not described here; you will find them in the description of the ODBC interface.

Table 38.49. sql_sizing Columns

Column Type	Description
sizing_id cardinal_number	Identifier of the sizing item
sizing_name character_data	Descriptive name of the sizing item
supported_value cardinal_number	Value of the sizing item, or 0 if the size is unlimited or cannot be determined, or null if the features for which the sizing item is applicable are not supported
comments character_data	Possibly a comment pertaining to the sizing item

38.52. table_constraints

The view `table_constraints` contains all constraints belonging to tables that the current user owns or has some privilege other than `SELECT` on.

Table 38.50. table_constraints Columns

Column Type	Description
constraint_catalog sql_identifier	Name of the database that contains the constraint (always the current database)
constraint_schema sql_identifier	Name of the schema that contains the constraint
constraint_name sql_identifier	Name of the constraint
table_catalog sql_identifier	Name of the database that contains the table (always the current database)
table_schema sql_identifier	Name of the schema that contains the table
table_name sql_identifier	Name of the table
constraint_type character_data	Type of the constraint: <code>CHECK</code> (includes not-null constraints), <code>FOREIGN KEY</code> , <code>PRIMARY KEY</code> , or <code>UNIQUE</code>
is_deferrable yes_or_no	<code>YES</code> if the constraint is deferrable, <code>NO</code> if not
initially_deferred yes_or_no	<code>YES</code> if the constraint is deferrable and initially deferred, <code>NO</code> if not
enforced yes_or_no	Applies to a feature not available in Postgres Pro (currently always <code>YES</code>)
nulls_distinct yes_or_no	If the constraint is a unique constraint, then <code>YES</code> if the constraint treats nulls as distinct or <code>NO</code> if it treats nulls as not distinct, otherwise null for other types of constraints.

38.53. table_privileges

The view `table_privileges` identifies all privileges granted on tables or views to a currently enabled role or by a currently enabled role. There is one row for each combination of table, grantor, and grantee.

Table 38.51. table_privileges Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to
table_catalog sql_identifier	Name of the database that contains the table (always the current database)
table_schema sql_identifier	Name of the schema that contains the table
table_name sql_identifier	Name of the table
privilege_type character_data	Type of the privilege: SELECT, INSERT, UPDATE, DELETE, TRUNCATE, REFERENCES, or TRIGGER
is_grantable yes_or_no	YES if the privilege is grantable, NO if not
with_hierarchy yes_or_no	In the SQL standard, WITH HIERARCHY OPTION is a separate (sub-)privilege allowing certain operations on table inheritance hierarchies. In Postgres Pro, this is included in the SELECT privilege, so this column shows YES if the privilege is SELECT, else NO.

38.54. tables

The view `tables` contains all tables and views defined in the current database. Only those tables and views are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.52. tables Columns

Column Type	Description
table_catalog sql_identifier	Name of the database that contains the table (always the current database)
table_schema sql_identifier	Name of the schema that contains the table
table_name sql_identifier	Name of the table
table_type character_data	Type of the table: BASE TABLE for a persistent base table (the normal table type), VIEW for a view, FOREIGN for a foreign table, or LOCAL TEMPORARY for a temporary table
self_referencing_column_name sql_identifier	Applies to a feature not available in Postgres Pro
reference_generation character_data	Applies to a feature not available in Postgres Pro
user_defined_type_catalog sql_identifier	If the table is a typed table, the name of the database that contains the underlying data type (always the current database), else null.
user_defined_type_schema sql_identifier	If the table is a typed table, the name of the schema that contains the underlying data type, else null.
user_defined_type_name sql_identifier	

Column Type	Description
	If the table is a typed table, the name of the underlying data type, else null.
is_insertable_into yes_or_no	YES if the table is insertable into, NO if not (Base tables are always insertable into, views not necessarily.)
is_typed yes_or_no	YES if the table is a typed table, NO if not
commit_action character_data	Not yet implemented

38.55. transforms

The view `transforms` contains information about the transforms defined in the current database. More precisely, it contains a row for each function contained in a transform (the “from SQL” or “to SQL” function).

Table 38.53. transforms Columns

Column Type	Description
udt_catalog sql_identifier	Name of the database that contains the type the transform is for (always the current database)
udt_schema sql_identifier	Name of the schema that contains the type the transform is for
udt_name sql_identifier	Name of the type the transform is for
specific_catalog sql_identifier	Name of the database containing the function (always the current database)
specific_schema sql_identifier	Name of the schema containing the function
specific_name sql_identifier	The “specific name” of the function. See Section 38.45 for more information.
group_name sql_identifier	The SQL standard allows defining transforms in “groups”, and selecting a group at run time. Postgres Pro does not support this. Instead, transforms are specific to a language. As a compromise, this field contains the language the transform is for.
transform_type character_data	FROM SQL or TO SQL

38.56. triggered_update_columns

For triggers in the current database that specify a column list (like `UPDATE OF column1, column2`), the view `triggered_update_columns` identifies these columns. Triggers that do not specify a column list are not included in this view. Only those columns are shown that the current user owns or has some privilege other than `SELECT` on.

Table 38.54. triggered_update_columns Columns

Column Type	Description
trigger_catalog sql_identifier	Name of the database that contains the trigger (always the current database)

Column Type	Description
trigger_schema sql_identifier	Name of the schema that contains the trigger
trigger_name sql_identifier	Name of the trigger
event_object_catalog sql_identifier	Name of the database that contains the table that the trigger is defined on (always the current database)
event_object_schema sql_identifier	Name of the schema that contains the table that the trigger is defined on
event_object_table sql_identifier	Name of the table that the trigger is defined on
event_object_column sql_identifier	Name of the column that the trigger is defined on

38.57. triggers

The view `triggers` contains all triggers defined in the current database on tables and views that the current user owns or has some privilege other than `SELECT` on.

Table 38.55. triggers Columns

Column Type	Description
trigger_catalog sql_identifier	Name of the database that contains the trigger (always the current database)
trigger_schema sql_identifier	Name of the schema that contains the trigger
trigger_name sql_identifier	Name of the trigger
event_manipulation character_data	Event that fires the trigger (INSERT, UPDATE, or DELETE)
event_object_catalog sql_identifier	Name of the database that contains the table that the trigger is defined on (always the current database)
event_object_schema sql_identifier	Name of the schema that contains the table that the trigger is defined on
event_object_table sql_identifier	Name of the table that the trigger is defined on
action_order cardinal_number	Firing order among triggers on the same table having the same <code>event_manipulation</code> , <code>action_timing</code> , and <code>action_orientation</code> . In Postgres Pro, triggers are fired in name order, so this column reflects that.
action_condition character_data	WHEN condition of the trigger, null if none (also null if the table is not owned by a currently enabled role)
action_statement character_data	Statement that is executed by the trigger (currently always <code>EXECUTE FUNCTION function(...)</code>)
action_orientation character_data	

Column Type	Description
	Identifies whether the trigger fires once for each processed row or once for each statement (ROW or STATEMENT)
action_timing character_data	Time at which the trigger fires (BEFORE, AFTER, or INSTEAD OF)
action_reference_old_table sql_identifier	Name of the “old” transition table, or null if none
action_reference_new_table sql_identifier	Name of the “new” transition table, or null if none
action_reference_old_row sql_identifier	Applies to a feature not available in Postgres Pro
action_reference_new_row sql_identifier	Applies to a feature not available in Postgres Pro
created_time_stamp	Applies to a feature not available in Postgres Pro

Triggers in Postgres Pro have two incompatibilities with the SQL standard that affect the representation in the information schema. First, trigger names are local to each table in Postgres Pro, rather than being independent schema objects. Therefore there can be duplicate trigger names defined in one schema, so long as they belong to different tables. (`trigger_catalog` and `trigger_schema` are really the values pertaining to the table that the trigger is defined on.) Second, triggers can be defined to fire on multiple events in Postgres Pro (e.g., ON INSERT OR UPDATE), whereas the SQL standard only allows one. If a trigger is defined to fire on multiple events, it is represented as multiple rows in the information schema, one for each type of event. As a consequence of these two issues, the primary key of the view `triggers` is really (`trigger_catalog`, `trigger_schema`, `event_object_table`, `trigger_name`, `event_manipulation`) instead of (`trigger_catalog`, `trigger_schema`, `trigger_name`), which is what the SQL standard specifies. Nonetheless, if you define your triggers in a manner that conforms with the SQL standard (trigger names unique in the schema and only one event type per trigger), this will not affect you.

Note

Prior to PostgreSQL 9.1, this view's columns `action_timing`, `action_reference_old_table`, `action_reference_new_table`, `action_reference_old_row`, and `action_reference_new_row` were named `condition_timing`, `condition_reference_old_table`, `condition_reference_new_table`, `condition_reference_old_row`, and `condition_reference_new_row` respectively. That was how they were named in the SQL:1999 standard. The new naming conforms to SQL:2003 and later.

38.58. udt_privileges

The view `udt_privileges` identifies `USAGE` privileges granted on user-defined types to a currently enabled role or by a currently enabled role. There is one row for each combination of type, grantor, and grantee. This view shows only composite types (see under [Section 38.60](#) for why); see [Section 38.59](#) for domain privileges.

Table 38.56. udt_privileges Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to

Column Type	Description
udt_catalog sql_identifier	Name of the database containing the type (always the current database)
udt_schema sql_identifier	Name of the schema containing the type
udt_name sql_identifier	Name of the type
privilege_type character_data	Always TYPE USAGE
is_grantable yes_or_no	YES if the privilege is grantable, NO if not

38.59. usage_privileges

The view `usage_privileges` identifies `USAGE` privileges granted on various kinds of objects to a currently enabled role or by a currently enabled role. In Postgres Pro, this currently applies to collations, domains, foreign-data wrappers, foreign servers, and sequences. There is one row for each combination of object, grantor, and grantee.

Since collations do not have real privileges in Postgres Pro, this view shows implicit non-grantable `USAGE` privileges granted by the owner to `PUBLIC` for all collations. The other object types, however, show real privileges.

In Postgres Pro, sequences also support `SELECT` and `UPDATE` privileges in addition to the `USAGE` privilege. These are nonstandard and therefore not visible in the information schema.

Table 38.57. usage_privileges Columns

Column Type	Description
grantor sql_identifier	Name of the role that granted the privilege
grantee sql_identifier	Name of the role that the privilege was granted to
object_catalog sql_identifier	Name of the database containing the object (always the current database)
object_schema sql_identifier	Name of the schema containing the object, if applicable, else an empty string
object_name sql_identifier	Name of the object
object_type character_data	COLLATION or DOMAIN or FOREIGN DATA WRAPPER or FOREIGN SERVER or SEQUENCE
privilege_type character_data	Always USAGE
is_grantable yes_or_no	YES if the privilege is grantable, NO if not

38.60. user_defined_types

The view `user_defined_types` currently contains all composite types defined in the current database. Only those types are shown that the current user has access to (by way of being the owner or having some privilege).

SQL knows about two kinds of user-defined types: structured types (also known as composite types in Postgres Pro) and distinct types (not implemented in Postgres Pro). To be future-proof, use the column `user_defined_type_category` to differentiate between these. Other user-defined types such as base types and enums, which are Postgres Pro extensions, are not shown here. For domains, see [Section 38.23](#) instead.

Table 38.58. `user_defined_types` Columns

Column Type		Description
<code>user_defined_type_catalog</code>	<code>sql_identifier</code>	Name of the database that contains the type (always the current database)
<code>user_defined_type_schema</code>	<code>sql_identifier</code>	Name of the schema that contains the type
<code>user_defined_type_name</code>	<code>sql_identifier</code>	Name of the type
<code>user_defined_type_category</code>	<code>character_data</code>	Currently always STRUCTURED
<code>is_instantiable</code>	<code>yes_or_no</code>	Applies to a feature not available in Postgres Pro
<code>is_final</code>	<code>yes_or_no</code>	Applies to a feature not available in Postgres Pro
<code>ordering_form</code>	<code>character_data</code>	Applies to a feature not available in Postgres Pro
<code>ordering_category</code>	<code>character_data</code>	Applies to a feature not available in Postgres Pro
<code>ordering_routine_catalog</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>ordering_routine_schema</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>ordering_routine_name</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>reference_type</code>	<code>character_data</code>	Applies to a feature not available in Postgres Pro
<code>data_type</code>	<code>character_data</code>	Applies to a feature not available in Postgres Pro
<code>character_maximum_length</code>	<code>cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>character_octet_length</code>	<code>cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>character_set_catalog</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_schema</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>character_set_name</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_catalog</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>collation_schema</code>	<code>sql_identifier</code>	Applies to a feature not available in Postgres Pro

Column Type	Description
<code>collation_name sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>numeric_precision cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>numeric_precision_radix cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>numeric_scale cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>datetime_precision cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>interval_type character_data</code>	Applies to a feature not available in Postgres Pro
<code>interval_precision cardinal_number</code>	Applies to a feature not available in Postgres Pro
<code>source_dtd_identifier sql_identifier</code>	Applies to a feature not available in Postgres Pro
<code>ref_dtd_identifier sql_identifier</code>	Applies to a feature not available in Postgres Pro

38.61. `user_mapping_options`

The view `user_mapping_options` contains all the options defined for user mappings in the current database. Only those user mappings are shown where the current user has access to the corresponding foreign server (by way of being the owner or having some privilege).

Table 38.59. `user_mapping_options` Columns

Column Type	Description
<code>authorization_identifier sql_identifier</code>	Name of the user being mapped, or <code>PUBLIC</code> if the mapping is public
<code>foreign_server_catalog sql_identifier</code>	Name of the database that the foreign server used by this mapping is defined in (always the current database)
<code>foreign_server_name sql_identifier</code>	Name of the foreign server used by this mapping
<code>option_name sql_identifier</code>	Name of an option
<code>option_value character_data</code>	Value of the option. This column will show as null unless the current user is the user being mapped, or the mapping is for <code>PUBLIC</code> and the current user is the server owner, or the current user is a superuser. The intent is to protect password information stored as user mapping option.

38.62. `user_mappings`

The view `user_mappings` contains all user mappings defined in the current database. Only those user mappings are shown where the current user has access to the corresponding foreign server (by way of being the owner or having some privilege).

Table 38.60. user_mappings Columns

Column Type	Description
authorization_identifier sql_identifier	Name of the user being mapped, or PUBLIC if the mapping is public
foreign_server_catalog sql_identifier	Name of the database that the foreign server used by this mapping is defined in (always the current database)
foreign_server_name sql_identifier	Name of the foreign server used by this mapping

38.63. view_column_usage

The view `view_column_usage` identifies all columns that are used in the query expression of a view (the `SELECT` statement that defines the view). A column is only included if the table that contains the column is owned by a currently enabled role.

Note

Columns of system tables are not included. This should be fixed sometime.

Table 38.61. view_column_usage Columns

Column Type	Description
view_catalog sql_identifier	Name of the database that contains the view (always the current database)
view_schema sql_identifier	Name of the schema that contains the view
view_name sql_identifier	Name of the view
table_catalog sql_identifier	Name of the database that contains the table that contains the column that is used by the view (always the current database)
table_schema sql_identifier	Name of the schema that contains the table that contains the column that is used by the view
table_name sql_identifier	Name of the table that contains the column that is used by the view
column_name sql_identifier	Name of the column that is used by the view

38.64. view_routine_usage

The view `view_routine_usage` identifies all routines (functions and procedures) that are used in the query expression of a view (the `SELECT` statement that defines the view). A routine is only included if that routine is owned by a currently enabled role.

Table 38.62. view_routine_usage Columns

Column Type	Description
table_catalog sql_identifier	

Column Type	Description
	Name of the database containing the view (always the current database)
table_schema sql_identifier	Name of the schema containing the view
table_name sql_identifier	Name of the view
specific_catalog sql_identifier	Name of the database containing the function (always the current database)
specific_schema sql_identifier	Name of the schema containing the function
specific_name sql_identifier	The “specific name” of the function. See Section 38.45 for more information.

38.65. view_table_usage

The view `view_table_usage` identifies all tables that are used in the query expression of a view (the `SELECT` statement that defines the view). A table is only included if that table is owned by a currently enabled role.

Note

System tables are not included. This should be fixed sometime.

Table 38.63. view_table_usage Columns

Column Type	Description
view_catalog sql_identifier	Name of the database that contains the view (always the current database)
view_schema sql_identifier	Name of the schema that contains the view
view_name sql_identifier	Name of the view
table_catalog sql_identifier	Name of the database that contains the table that is used by the view (always the current database)
table_schema sql_identifier	Name of the schema that contains the table that is used by the view
table_name sql_identifier	Name of the table that is used by the view

38.66. views

The view `views` contains all views defined in the current database. Only those views are shown that the current user has access to (by way of being the owner or having some privilege).

Table 38.64. views Columns

Column Type	Description
table_catalog sql_identifier	

Column Type	Description
	Name of the database that contains the view (always the current database)
table_schema sql_identifier	Name of the schema that contains the view
table_name sql_identifier	Name of the view
view_definition character_data	Query expression defining the view (null if the view is not owned by a currently enabled role)
check_option character_data	CASCADED or LOCAL if the view has a CHECK OPTION defined on it, NONE if not
is_updatable yes_or_no	YES if the view is updatable (allows UPDATE and DELETE), NO if not
is_insertable_into yes_or_no	YES if the view is insertable into (allows INSERT), NO if not
is_trigger_updatable yes_or_no	YES if the view has an INSTEAD OF UPDATE trigger defined on it, NO if not
is_trigger_deletable yes_or_no	YES if the view has an INSTEAD OF DELETE trigger defined on it, NO if not
is_trigger_insertable_into yes_or_no	YES if the view has an INSTEAD OF INSERT trigger defined on it, NO if not

Part V. Server Programming

This part is about extending the server functionality with user-defined functions, data types, triggers, etc. These are advanced topics which should be approached only after all the other user documentation about Postgres Pro has been understood. Later chapters in this part describe the server-side programming languages available in the Postgres Pro distribution as well as general issues concerning server-side programming. It is essential to read at least the earlier sections of [Chapter 39](#) (covering functions) before diving into the material about server-side programming.

Chapter 39. Extending SQL

In the sections that follow, we will discuss how you can extend the Postgres Pro SQL query language by adding:

- functions (starting in [Section 39.3](#))
- aggregates (starting in [Section 39.12](#))
- data types (starting in [Section 39.13](#))
- operators (starting in [Section 39.14](#))
- operator classes for indexes (starting in [Section 39.16](#))
- packages of related objects (starting in [Section 39.17](#))

39.1. How Extensibility Works

Postgres Pro is extensible because its operation is catalog-driven. If you are familiar with standard relational database systems, you know that they store information about databases, tables, columns, etc., in what are commonly known as system catalogs. (Some systems call this the data dictionary.) The catalogs appear to the user as tables like any other, but the DBMS stores its internal bookkeeping in them. One key difference between Postgres Pro and standard relational database systems is that Postgres Pro stores much more information in its catalogs: not only information about tables and columns, but also information about data types, functions, access methods, and so on. These tables can be modified by the user, and since Postgres Pro bases its operation on these tables, this means that Postgres Pro can be extended by users. By comparison, conventional database systems can only be extended by changing hardcoded procedures in the source code or by loading modules specially written by the DBMS vendor.

The Postgres Pro server can moreover incorporate user-written code into itself through dynamic loading. That is, the user can specify an object code file (e.g., a shared library) that implements a new type or function, and Postgres Pro will load it as required. Code written in SQL is even more trivial to add to the server. This ability to modify its operation “on the fly” makes Postgres Pro uniquely suited for rapid prototyping of new applications and storage structures.

39.2. The Postgres Pro Type System

Postgres Pro data types can be divided into base types, container types, domains, and pseudo-types.

39.2.1. Base Types

Base types are those, like `integer`, that are implemented below the level of the SQL language (typically in a low-level language such as C). They generally correspond to what are often known as abstract data types. Postgres Pro can only operate on such types through functions provided by the user and only understands the behavior of such types to the extent that the user describes them. The built-in base types are described in [Chapter 8](#).

Enumerated (enum) types can be considered as a subcategory of base types. The main difference is that they can be created using just SQL commands, without any low-level programming. Refer to [Section 8.7](#) for more information.

39.2.2. Container Types

Postgres Pro has three kinds of “container” types, which are types that contain multiple values of other types. These are arrays, composites, and ranges.

Arrays can hold multiple values that are all of the same type. An array type is automatically created for each base type, composite type, range type, and domain type. But there are no arrays of arrays. So far as the type system is concerned, multi-dimensional arrays are the same as one-dimensional arrays. Refer to [Section 8.15](#) for more information.

Composite types, or row types, are created whenever the user creates a table. It is also possible to use `CREATE TYPE` to define a “stand-alone” composite type with no associated table. A composite type is

simply a list of types with associated field names. A value of a composite type is a row or record of field values. Refer to [Section 8.16](#) for more information.

A range type can hold two values of the same type, which are the lower and upper bounds of the range. Range types are user-created, although a few built-in ones exist. Refer to [Section 8.17](#) for more information.

39.2.3. Domains

A domain is based on a particular underlying type and for many purposes is interchangeable with its underlying type. However, a domain can have constraints that restrict its valid values to a subset of what the underlying type would allow. Domains are created using the SQL command [CREATE DOMAIN](#). Refer to [Section 8.18](#) for more information.

39.2.4. Pseudo-Types

There are a few “pseudo-types” for special purposes. Pseudo-types cannot appear as columns of tables or components of container types, but they can be used to declare the argument and result types of functions. This provides a mechanism within the type system to identify special classes of functions. [Table 8.27](#) lists the existing pseudo-types.

39.2.5. Polymorphic Types

Some pseudo-types of special interest are the *polymorphic types*, which are used to declare *polymorphic functions*. This powerful feature allows a single function definition to operate on many different data types, with the specific data type(s) being determined by the data types actually passed to it in a particular call. The polymorphic types are shown in [Table 39.1](#). Some examples of their use appear in [Section 39.5.11](#).

Table 39.1. Polymorphic Types

Name	Family	Description
anyelement	Simple	Indicates that a function accepts any data type
anyarray	Simple	Indicates that a function accepts any array data type
anynonarray	Simple	Indicates that a function accepts any non-array data type
anyenum	Simple	Indicates that a function accepts any enum data type (see Section 8.7)
anyrange	Simple	Indicates that a function accepts any range data type (see Section 8.17)
anymultirange	Simple	Indicates that a function accepts any multirange data type (see Section 8.17)
anycompatible	Common	Indicates that a function accepts any data type, with automatic promotion of multiple arguments to a common data type
anycompatiblearray	Common	Indicates that a function accepts any array data type, with automatic promotion of multiple arguments to a common data type
anycompatiblenonarray	Common	Indicates that a function accepts any non-array data type, with automatic promotion of multiple arguments to a common data type

Name	Family	Description
<code>anycompatiblerange</code>	Common	Indicates that a function accepts any range data type, with automatic promotion of multiple arguments to a common data type
<code>anycompatiblemultirange</code>	Common	Indicates that a function accepts any multirange data type, with automatic promotion of multiple arguments to a common data type

Polymorphic arguments and results are tied to each other and are resolved to specific data types when a query calling a polymorphic function is parsed. When there is more than one polymorphic argument, the actual data types of the input values must match up as described below. If the function's result type is polymorphic, or it has output parameters of polymorphic types, the types of those results are deduced from the actual types of the polymorphic inputs as described below.

For the “simple” family of polymorphic types, the matching and deduction rules work like this:

Each position (either argument or return value) declared as `anyelement` is allowed to have any specific actual data type, but in any given call they must all be the *same* actual type. Each position declared as `anyarray` can have any array data type, but similarly they must all be the same type. And similarly, positions declared as `anyrange` must all be the same range type. Likewise for `anymultirange`.

Furthermore, if there are positions declared `anyarray` and others declared `anyelement`, the actual array type in the `anyarray` positions must be an array whose elements are the same type appearing in the `anyelement` positions. `anynonarray` is treated exactly the same as `anyelement`, but adds the additional constraint that the actual type must not be an array type. `anyenum` is treated exactly the same as `anyelement`, but adds the additional constraint that the actual type must be an enum type.

Similarly, if there are positions declared `anyrange` and others declared `anyelement` or `anyarray`, the actual range type in the `anyrange` positions must be a range whose subtype is the same type appearing in the `anyelement` positions and the same as the element type of the `anyarray` positions. If there are positions declared `anymultirange`, their actual multirange type must contain ranges matching parameters declared `anyrange` and base elements matching parameters declared `anyelement` and `anyarray`.

Thus, when more than one argument position is declared with a polymorphic type, the net effect is that only certain combinations of actual argument types are allowed. For example, a function declared as `equal(anyelement, anyelement)` will take any two input values, so long as they are of the same data type.

When the return value of a function is declared as a polymorphic type, there must be at least one argument position that is also polymorphic, and the actual data type(s) supplied for the polymorphic arguments determine the actual result type for that call. For example, if there were not already an array subscripting mechanism, one could define a function that implements subscripting as `subscript(anyarray, integer) returns anyelement`. This declaration constrains the actual first argument to be an array type, and allows the parser to infer the correct result type from the actual first argument's type. Another example is that a function declared as `f(anyarray) returns anyenum` will only accept arrays of enum types.

In most cases, the parser can infer the actual data type for a polymorphic result type from arguments that are of a different polymorphic type in the same family; for example `anyarray` can be deduced from `anyelement` or vice versa. An exception is that a polymorphic result of type `anyrange` requires an argument of type `anyrange`; it cannot be deduced from `anyarray` or `anyelement` arguments. This is because there could be multiple range types with the same subtype.

Note that `anynonarray` and `anyenum` do not represent separate type variables; they are the same type as `anyelement`, just with an additional constraint. For example, declaring a function as `f(anyelement,`

`anyenum`) is equivalent to declaring it as `f(anyenum, anyenum)`: both actual arguments have to be the same enum type.

For the “common” family of polymorphic types, the matching and deduction rules work approximately the same as for the “simple” family, with one major difference: the actual types of the arguments need not be identical, so long as they can be implicitly cast to a single common type. The common type is selected following the same rules as for `UNION` and related constructs (see [Section 10.5](#)). Selection of the common type considers the actual types of `anycompatible` and `anycompatiblenonarray` inputs, the array element types of `anycompatiblearray` inputs, the range subtypes of `anycompatiblerange` inputs, and the multirange subtypes of `anycompatiblemultirange` inputs. If `anycompatiblenonarray` is present then the common type is required to be a non-array type. Once a common type is identified, arguments in `anycompatible` and `anycompatiblenonarray` positions are automatically cast to that type, and arguments in `anycompatiblearray` positions are automatically cast to the array type for that type.

Since there is no way to select a range type knowing only its subtype, use of `anycompatiblerange` and/or `anycompatiblemultirange` requires that all arguments declared with that type have the same actual range and/or multirange type, and that that type's subtype agree with the selected common type, so that no casting of the range values is required. As with `anyrange` and `anymultirange`, use of `anycompatiblerange` and `anymultirange` as a function result type requires that there be an `anycompatiblerange` or `anycompatiblemultirange` argument.

Notice that there is no `anycompatibleenum` type. Such a type would not be very useful, since there normally are not any implicit casts to enum types, meaning that there would be no way to resolve a common type for dissimilar enum inputs.

The “simple” and “common” polymorphic families represent two independent sets of type variables. Consider for example

```
CREATE FUNCTION myfunc(a anyelement, b anyelement,
                      c anycompatible, d anycompatible)
RETURNS anycompatible AS ...
```

In an actual call of this function, the first two inputs must have exactly the same type. The last two inputs must be promotable to a common type, but this type need not have anything to do with the type of the first two inputs. The result will have the common type of the last two inputs.

A variadic function (one taking a variable number of arguments, as in [Section 39.5.6](#)) can be polymorphic: this is accomplished by declaring its last parameter as `VARIADIC anyarray` or `VARIADIC anycompatiblearray`. For purposes of argument matching and determining the actual result type, such a function behaves the same as if you had written the appropriate number of `anynonarray` or `anycompatiblenonarray` parameters.

39.3. User-Defined Functions

Postgres Pro provides four kinds of functions:

- query language functions (functions written in SQL) ([Section 39.5](#))
- procedural language functions (functions written in, for example, PL/pgSQL or PL/Tcl) ([Section 39.8](#))
- internal functions ([Section 39.9](#))
- C-language functions ([Section 39.10](#))

Every kind of function can take base types, composite types, or combinations of these as arguments (parameters). In addition, every kind of function can return a base type or a composite type. Functions can also be defined to return sets of base or composite values.

Many kinds of functions can take or return certain pseudo-types (such as polymorphic types), but the available facilities vary. Consult the description of each kind of function for more details.

It's easiest to define SQL functions, so we'll start by discussing those. Most of the concepts presented for SQL functions will carry over to the other types of functions.

Throughout this chapter, it can be useful to look at the reference page of the [CREATE FUNCTION](#) command to understand the examples better.

39.4. User-Defined Procedures

A procedure is a database object similar to a function. The key differences are:

- Procedures are defined with the [CREATE PROCEDURE](#) command, not [CREATE FUNCTION](#).
- Procedures do not return a function value; hence [CREATE PROCEDURE](#) lacks a [RETURNS](#) clause. However, procedures can instead return data to their callers via output parameters.
- While a function is called as part of a query or DML command, a procedure is called in isolation using the [CALL](#) command.
- A procedure can commit or roll back transactions during its execution (then automatically beginning a new transaction), so long as the invoking [CALL](#) command is not part of an explicit transaction block. A function cannot do that.
- Certain function attributes, such as strictness, don't apply to procedures. Those attributes control how the function is used in a query, which isn't relevant to procedures.

The explanations in the following sections about how to define user-defined functions apply to procedures as well, except for the points made above.

Collectively, functions and procedures are also known as *routines*. There are commands such as [ALTER ROUTINE](#) and [DROP ROUTINE](#) that can operate on functions and procedures without having to know which kind it is. Note, however, that there is no [CREATE ROUTINE](#) command.

39.5. Query Language (SQL) Functions

SQL functions execute an arbitrary list of SQL statements, returning the result of the last query in the list. In the simple (non-set) case, the first row of the last query's result will be returned. (Bear in mind that “the first row” of a multirow result is not well-defined unless you use [ORDER BY](#).) If the last query happens to return no rows at all, the null value will be returned.

Alternatively, an SQL function can be declared to return a set (that is, multiple rows) by specifying the function's return type as [SETOF sometype](#), or equivalently by declaring it as [RETURNS TABLE\(columns\)](#). In this case all rows of the last query's result are returned. Further details appear below.

The body of an SQL function must be a list of SQL statements separated by semicolons. A semicolon after the last statement is optional. Unless the function is declared to return [void](#), the last statement must be a [SELECT](#), or an [INSERT](#), [UPDATE](#), [DELETE](#), or [MERGE](#) that has a [RETURNING](#) clause.

Any collection of commands in the SQL language can be packaged together and defined as a function. Besides [SELECT](#) queries, the commands can include data modification queries ([INSERT](#), [UPDATE](#), [DELETE](#), and [MERGE](#)), as well as other SQL commands. (You cannot use transaction control commands, e.g., [COMMIT](#), [SAVEPOINT](#), and some utility commands, e.g., [VACUUM](#), in SQL functions.) However, the final command must be a [SELECT](#) or have a [RETURNING](#) clause that returns whatever is specified as the function's return type. Alternatively, if you want to define an SQL function that performs actions but has no useful value to return, you can define it as returning [void](#). For example, this function removes rows with negative salaries from the `emp` table:

```
CREATE FUNCTION clean_emp() RETURNS void AS '
    DELETE FROM emp
        WHERE salary < 0;
' LANGUAGE SQL;

SELECT clean_emp();
```

```
clean_emp
-----

(1 row)
```

You can also write this as a procedure, thus avoiding the issue of the return type. For example:

```
CREATE PROCEDURE clean_emp() AS '
    DELETE FROM emp
        WHERE salary < 0;
' LANGUAGE SQL;

CALL clean_emp();
```

In simple cases like this, the difference between a function returning `void` and a procedure is mostly stylistic. However, procedures offer additional functionality such as transaction control that is not available in functions. Also, procedures are SQL standard whereas returning `void` is a Postgres Pro extension.

Note

The entire body of an SQL function is parsed before any of it is executed. While an SQL function can contain commands that alter the system catalogs (e.g., `CREATE TABLE`), the effects of such commands will not be visible during parse analysis of later commands in the function. Thus, for example, `CREATE TABLE foo (...); INSERT INTO foo VALUES (...);` will not work as desired if packaged up into a single SQL function, since `foo` won't exist yet when the `INSERT` command is parsed. It's recommended to use PL/pgSQL instead of an SQL function in this type of situation.

The syntax of the `CREATE FUNCTION` command requires the function body to be written as a string constant. It is usually most convenient to use dollar quoting (see [Section 4.1.2.4](#)) for the string constant. If you choose to use regular single-quoted string constant syntax, you must double single quote marks (') and backslashes (\) (assuming escape string syntax) in the body of the function (see [Section 4.1.2.1](#)).

39.5.1. Arguments for SQL Functions

Arguments of an SQL function can be referenced in the function body using either names or numbers. Examples of both methods appear below.

To use a name, declare the function argument as having a name, and then just write that name in the function body. If the argument name is the same as any column name in the current SQL command within the function, the column name will take precedence. To override this, qualify the argument name with the name of the function itself, that is `function_name.argument_name`. (If this would conflict with a qualified column name, again the column name wins. You can avoid the ambiguity by choosing a different alias for the table within the SQL command.)

In the older numeric approach, arguments are referenced using the syntax `$n`: `$1` refers to the first input argument, `$2` to the second, and so on. This will work whether or not the particular argument was declared with a name.

If an argument is of a composite type, then the dot notation, e.g., `argname.fieldname` or `$1.fieldname`, can be used to access attributes of the argument. Again, you might need to qualify the argument's name with the function name to make the form with an argument name unambiguous.

SQL function arguments can only be used as data values, not as identifiers. Thus for example this is reasonable:

```
INSERT INTO mytable VALUES ($1);
```

but this will not work:

```
INSERT INTO $1 VALUES (42);
```

Note

The ability to use names to reference SQL function arguments was added in PostgreSQL 9.2. Functions to be used in older servers must use the $\$n$ notation.

39.5.2. SQL Functions on Base Types

The simplest possible SQL function has no arguments and simply returns a base type, such as `integer`:

```
CREATE FUNCTION one() RETURNS integer AS $$
    SELECT 1 AS result;
$$ LANGUAGE SQL;
```

```
-- Alternative syntax for string literal:
CREATE FUNCTION one() RETURNS integer AS '
    SELECT 1 AS result;
' LANGUAGE SQL;
```

```
SELECT one();
```

```
one
-----
1
```

Notice that we defined a column alias within the function body for the result of the function (with the name `result`), but this column alias is not visible outside the function. Hence, the result is labeled `one` instead of `result`.

It is almost as easy to define SQL functions that take base types as arguments:

```
CREATE FUNCTION add_em(x integer, y integer) RETURNS integer AS $$
    SELECT x + y;
$$ LANGUAGE SQL;
```

```
SELECT add_em(1, 2) AS answer;
```

```
answer
-----
3
```

Alternatively, we could dispense with names for the arguments and use numbers:

```
CREATE FUNCTION add_em(integer, integer) RETURNS integer AS $$
    SELECT $1 + $2;
$$ LANGUAGE SQL;
```

```
SELECT add_em(1, 2) AS answer;
```

```
answer
-----
3
```

Here is a more useful function, which might be used to debit a bank account:

```
CREATE FUNCTION tf1 (accountno integer, debit numeric) RETURNS numeric AS $$
    UPDATE bank
        SET balance = balance - debit
        WHERE accountno = tf1.accountno;
    SELECT 1;
$$ LANGUAGE SQL;
```

A user could execute this function to debit account 17 by \$100.00 as follows:

```
SELECT tf1(17, 100.0);
```

In this example, we chose the name `accountno` for the first argument, but this is the same as the name of a column in the `bank` table. Within the `UPDATE` command, `accountno` refers to the column `bank.accountno`, so `tf1.accountno` must be used to refer to the argument. We could of course avoid this by using a different name for the argument.

In practice one would probably like a more useful result from the function than a constant 1, so a more likely definition is:

```
CREATE FUNCTION tf1 (accountno integer, debit numeric) RETURNS numeric AS $$
    UPDATE bank
        SET balance = balance - debit
        WHERE accountno = tf1.accountno;
    SELECT balance FROM bank WHERE accountno = tf1.accountno;
$$ LANGUAGE SQL;
```

which adjusts the balance and returns the new balance. The same thing could be done in one command using `RETURNING`:

```
CREATE FUNCTION tf1 (accountno integer, debit numeric) RETURNS numeric AS $$
    UPDATE bank
        SET balance = balance - debit
        WHERE accountno = tf1.accountno
    RETURNING balance;
$$ LANGUAGE SQL;
```

If the final `SELECT` or `RETURNING` clause in an SQL function does not return exactly the function's declared result type, Postgres Pro will automatically cast the value to the required type, if that is possible with an implicit or assignment cast. Otherwise, you must write an explicit cast. For example, suppose we wanted the previous `add_em` function to return type `float8` instead. It's sufficient to write

```
CREATE FUNCTION add_em(integer, integer) RETURNS float8 AS $$
    SELECT $1 + $2;
$$ LANGUAGE SQL;
```

since the `integer` sum can be implicitly cast to `float8`. (See [Chapter 10](#) or [CREATE CAST](#) for more about casts.)

39.5.3. SQL Functions on Composite Types

When writing functions with arguments of composite types, we must not only specify which argument we want but also the desired attribute (field) of that argument. For example, suppose that `emp` is a table containing employee data, and therefore also the name of the composite type of each row of the table. Here is a function `double_salary` that computes what someone's salary would be if it were doubled:

```
CREATE TABLE emp (
    name      text,
    salary     numeric,
    age        integer,
    cubicle    point
);

INSERT INTO emp VALUES ('Bill', 4200, 45, '(2,1)');

CREATE FUNCTION double_salary(emp) RETURNS numeric AS $$
    SELECT $1.salary * 2 AS salary;
$$ LANGUAGE SQL;

SELECT name, double_salary(emp.*) AS dream
```



```
FROM emp
WHERE emp.cubicle ~= point '(2,1)';
```

```
name | dream
-----+-----
Bill | 8400
```

Notice the use of the syntax `$1.salary` to select one field of the argument row value. Also notice how the calling `SELECT` command uses `table_name.*` to select the entire current row of a table as a composite value. The table row can alternatively be referenced using just the table name, like this:

```
SELECT name, double_salary(emp) AS dream
FROM emp
WHERE emp.cubicle ~= point '(2,1)';
```

but this usage is deprecated since it's easy to get confused. (See [Section 8.16.5](#) for details about these two notations for the composite value of a table row.)

Sometimes it is handy to construct a composite argument value on-the-fly. This can be done with the `ROW` construct. For example, we could adjust the data being passed to the function:

```
SELECT name, double_salary(ROW(name, salary*1.1, age, cubicle)) AS dream
FROM emp;
```

It is also possible to build a function that returns a composite type. This is an example of a function that returns a single `emp` row:

```
CREATE FUNCTION new_emp() RETURNS emp AS $$
SELECT text 'None' AS name,
       1000.0 AS salary,
       25 AS age,
       point '(2,2)' AS cubicle;
$$ LANGUAGE SQL;
```

In this example we have specified each of the attributes with a constant value, but any computation could have been substituted for these constants.

Note two important things about defining the function:

- The select list order in the query must be exactly the same as that in which the columns appear in the composite type. (Naming the columns, as we did above, is irrelevant to the system.)
- We must ensure each expression's type can be cast to that of the corresponding column of the composite type. Otherwise we'll get errors like this:

```
ERROR:  return type mismatch in function declared to return emp
DETAIL:  Final statement returns text instead of point at column 4.
```

As with the base-type case, the system will not insert explicit casts automatically, only implicit or assignment casts.

A different way to define the same function is:

```
CREATE FUNCTION new_emp() RETURNS emp AS $$
SELECT ROW('None', 1000.0, 25, '(2,2)')::emp;
$$ LANGUAGE SQL;
```

Here we wrote a `SELECT` that returns just a single column of the correct composite type. This isn't really better in this situation, but it is a handy alternative in some cases — for example, if we need to compute the result by calling another function that returns the desired composite value. Another example is that if we are trying to write a function that returns a domain over composite, rather than a plain composite type, it is always necessary to write it as returning a single column, since there is no way to cause a coercion of the whole row result.

We could call this function directly either by using it in a value expression:

```
SELECT new_emp();
```

```

      new_emp
-----
(None,1000.0,25,"(2,2)")

```

or by calling it as a table function:

```
SELECT * FROM new_emp();
```

```

name | salary | age | cubicle
-----+-----+-----+-----
None | 1000.0 | 25  | (2,2)

```

The second way is described more fully in [Section 39.5.8](#).

When you use a function that returns a composite type, you might want only one field (attribute) from its result. You can do that with syntax like this:

```
SELECT (new_emp()).name;
```

```

name
-----
None

```

The extra parentheses are needed to keep the parser from getting confused. If you try to do it without them, you get something like this:

```

SELECT new_emp().name;
ERROR:  syntax error at or near "."
LINE 1: SELECT new_emp().name;
                        ^

```

Another option is to use functional notation for extracting an attribute:

```
SELECT name(new_emp());
```

```

name
-----
None

```

As explained in [Section 8.16.5](#), the field notation and functional notation are equivalent.

Another way to use a function returning a composite type is to pass the result to another function that accepts the correct row type as input:

```

CREATE FUNCTION getname(emp) RETURNS text AS $$
    SELECT $1.name;
$$ LANGUAGE SQL;

```

```
SELECT getname(new_emp());
```

```

getname
-----
None
(1 row)

```

39.5.4. SQL Functions with Output Parameters

An alternative way of describing a function's results is to define it with *output parameters*, as in this example:

```

CREATE FUNCTION add_em (IN x int, IN y int, OUT sum int)
AS 'SELECT x + y'
LANGUAGE SQL;

```

```
SELECT add_em(3,7);
      add_em
-----
         10
(1 row)
```

This is not essentially different from the version of `add_em` shown in [Section 39.5.2](#). The real value of output parameters is that they provide a convenient way of defining functions that return several columns. For example,

```
CREATE FUNCTION sum_n_product (x int, y int, OUT sum int, OUT product int)
AS 'SELECT x + y, x * y'
LANGUAGE SQL;
```

```
SELECT * FROM sum_n_product(11,42);
      sum | product
-----+-----
        53 |      462
(1 row)
```

What has essentially happened here is that we have created an anonymous composite type for the result of the function. The above example has the same end result as

```
CREATE TYPE sum_prod AS (sum int, product int);

CREATE FUNCTION sum_n_product (int, int) RETURNS sum_prod
AS 'SELECT $1 + $2, $1 * $2'
LANGUAGE SQL;
```

but not having to bother with the separate composite type definition is often handy. Notice that the names attached to the output parameters are not just decoration, but determine the column names of the anonymous composite type. (If you omit a name for an output parameter, the system will choose a name on its own.)

Notice that output parameters are not included in the calling argument list when invoking such a function from SQL. This is because Postgres Pro considers only the input parameters to define the function's calling signature. That means also that only the input parameters matter when referencing the function for purposes such as dropping it. We could drop the above function with either of

```
DROP FUNCTION sum_n_product (x int, y int, OUT sum int, OUT product int);
DROP FUNCTION sum_n_product (int, int);
```

Parameters can be marked as `IN` (the default), `OUT`, `INOUT`, or `VARIADIC`. An `INOUT` parameter serves as both an input parameter (part of the calling argument list) and an output parameter (part of the result record type). `VARIADIC` parameters are input parameters, but are treated specially as described below.

39.5.5. SQL Procedures with Output Parameters

Output parameters are also supported in procedures, but they work a bit differently from functions. In `CALL` commands, output parameters must be included in the argument list. For example, the bank account debiting routine from earlier could be written like this:

```
CREATE PROCEDURE tp1 (accountno integer, debit numeric, OUT new_balance numeric) AS $$
    UPDATE bank
        SET balance = balance - debit
        WHERE accountno = tp1.accountno
    RETURNING balance;
$$ LANGUAGE SQL;
```

To call this procedure, an argument matching the `OUT` parameter must be included. It's customary to write `NULL`:

```
CALL tp1(17, 100.0, NULL);
```

If you write something else, it must be an expression that is implicitly coercible to the declared type of the parameter, just as for input parameters. Note however that such an expression will not be evaluated.

When calling a procedure from PL/pgSQL, instead of writing `NULL` you must write a variable that will receive the procedure's output. See [Section 44.6.3](#) for details.

39.5.6. SQL Functions with Variable Numbers of Arguments

SQL functions can be declared to accept variable numbers of arguments, so long as all the “optional” arguments are of the same data type. The optional arguments will be passed to the function as an array. The function is declared by marking the last parameter as `VARIADIC`; this parameter must be declared as being of an array type. For example:

```
CREATE FUNCTION mleast(VARIADIC arr numeric[]) RETURNS numeric AS $$
    SELECT min($1[i]) FROM generate_subscripts($1, 1) g(i);
$$ LANGUAGE SQL;
```

```
SELECT mleast(10, -1, 5, 4.4);
 mleast
-----
      -1
(1 row)
```

Effectively, all the actual arguments at or beyond the `VARIADIC` position are gathered up into a one-dimensional array, as if you had written

```
SELECT mleast(ARRAY[10, -1, 5, 4.4]);    -- doesn't work
```

You can't actually write that, though — or at least, it will not match this function definition. A parameter marked `VARIADIC` matches one or more occurrences of its element type, not of its own type.

Sometimes it is useful to be able to pass an already-constructed array to a variadic function; this is particularly handy when one variadic function wants to pass on its array parameter to another one. Also, this is the only secure way to call a variadic function found in a schema that permits untrusted users to create objects; see [Section 10.3](#). You can do this by specifying `VARIADIC` in the call:

```
SELECT mleast(VARIADIC ARRAY[10, -1, 5, 4.4]);
```

This prevents expansion of the function's variadic parameter into its element type, thereby allowing the array argument value to match normally. `VARIADIC` can only be attached to the last actual argument of a function call.

Specifying `VARIADIC` in the call is also the only way to pass an empty array to a variadic function, for example:

```
SELECT mleast(VARIADIC ARRAY[]::numeric[]);
```

Simply writing `SELECT mleast()` does not work because a variadic parameter must match at least one actual argument. (You could define a second function also named `mleast`, with no parameters, if you wanted to allow such calls.)

The array element parameters generated from a variadic parameter are treated as not having any names of their own. This means it is not possible to call a variadic function using named arguments ([Section 4.3](#)), except when you specify `VARIADIC`. For example, this will work:

```
SELECT mleast(VARIADIC arr => ARRAY[10, -1, 5, 4.4]);
```

but not these:

```
SELECT mleast(arr => 10);
SELECT mleast(arr => ARRAY[10, -1, 5, 4.4]);
```

39.5.7. SQL Functions with Default Values for Arguments

Functions can be declared with default values for some or all input arguments. The default values are inserted whenever the function is called with insufficiently many actual arguments. Since arguments can only be omitted from the end of the actual argument list, all parameters after a parameter with a default

value have to have default values as well. (Although the use of named argument notation could allow this restriction to be relaxed, it's still enforced so that positional argument notation works sensibly.) Whether or not you use it, this capability creates a need for precautions when calling functions in databases where some users mistrust other users; see [Section 10.3](#).

For example:

```
CREATE FUNCTION foo(a int, b int DEFAULT 2, c int DEFAULT 3)
RETURNS int
LANGUAGE SQL
AS $$
    SELECT $1 + $2 + $3;
$$;

SELECT foo(10, 20, 30);
foo
-----
60
(1 row)

SELECT foo(10, 20);
foo
-----
33
(1 row)

SELECT foo(10);
foo
-----
15
(1 row)

SELECT foo(); -- fails since there is no default for the first argument
ERROR:  function foo() does not exist
```

The = sign can also be used in place of the key word `DEFAULT`.

39.5.8. SQL Functions as Table Sources

All SQL functions can be used in the `FROM` clause of a query, but it is particularly useful for functions returning composite types. If the function is defined to return a base type, the table function produces a one-column table. If the function is defined to return a composite type, the table function produces a column for each attribute of the composite type.

Here is an example:

```
CREATE TABLE foo (fooid int, foosubid int, foename text);
INSERT INTO foo VALUES (1, 1, 'Joe');
INSERT INTO foo VALUES (1, 2, 'Ed');
INSERT INTO foo VALUES (2, 1, 'Mary');

CREATE FUNCTION getfoo(int) RETURNS foo AS $$
    SELECT * FROM foo WHERE fooid = $1;
$$ LANGUAGE SQL;

SELECT *, upper(foename) FROM getfoo(1) AS t1;

fooid | foosubid | foename | upper
-----+-----+-----+-----
1 | 1 | Joe | JOE
(1 row)
```

As the example shows, we can work with the columns of the function's result just the same as if they were columns of a regular table.

Note that we only got one row out of the function. This is because we did not use `SETOF`. That is described in the next section.

39.5.9. SQL Functions Returning Sets

When an SQL function is declared as returning `SETOF sometype`, the function's final query is executed to completion, and each row it outputs is returned as an element of the result set.

This feature is normally used when calling the function in the `FROM` clause. In this case each row returned by the function becomes a row of the table seen by the query. For example, assume that table `foo` has the same contents as above, and we say:

```
CREATE FUNCTION getfoo(int) RETURNS SETOF foo AS $$
    SELECT * FROM foo WHERE fooid = $1;
$$ LANGUAGE SQL;
```

```
SELECT * FROM getfoo(1) AS t1;
```

Then we would get:

fooid	foosubid	foaname
1	1	Joe
1	2	Ed

(2 rows)

It is also possible to return multiple rows with the columns defined by output parameters, like this:

```
CREATE TABLE tab (y int, z int);
INSERT INTO tab VALUES (1, 2), (3, 4), (5, 6), (7, 8);

CREATE FUNCTION sum_n_product_with_tab (x int, OUT sum int, OUT product int)
RETURNS SETOF record
AS $$
    SELECT $1 + tab.y, $1 * tab.y FROM tab;
$$ LANGUAGE SQL;
```

```
SELECT * FROM sum_n_product_with_tab(10);
sum | product
-----+-----
11 | 10
13 | 30
15 | 50
17 | 70
(4 rows)
```

The key point here is that you must write `RETURNS SETOF record` to indicate that the function returns multiple rows instead of just one. If there is only one output parameter, write that parameter's type instead of `record`.

It is frequently useful to construct a query's result by invoking a set-returning function multiple times, with the parameters for each invocation coming from successive rows of a table or subquery. The preferred way to do this is to use the `LATERAL` key word, which is described in [Section 7.2.1.5](#). Here is an example using a set-returning function to enumerate elements of a tree structure:

```
SELECT * FROM nodes;
name | parent
-----+-----
Top |
Child1 | Top
```

```
Child2      | Top
Child3      | Top
SubChild1   | Child1
SubChild2   | Child1
(6 rows)
```

```
CREATE FUNCTION listchildren(text) RETURNS SETOF text AS $$
    SELECT name FROM nodes WHERE parent = $1
$$ LANGUAGE SQL STABLE;
```

```
SELECT * FROM listchildren('Top');
listchildren
```

```
-----
Child1
Child2
Child3
(3 rows)
```

```
SELECT name, child FROM nodes, LATERAL listchildren(name) AS child;
```

```
name | child
-----+-----
Top   | Child1
Top   | Child2
Top   | Child3
Child1 | SubChild1
Child1 | SubChild2
(5 rows)
```

This example does not do anything that we couldn't have done with a simple join, but in more complex calculations the option to put some of the work into a function can be quite convenient.

Functions returning sets can also be called in the select list of a query. For each row that the query generates by itself, the set-returning function is invoked, and an output row is generated for each element of the function's result set. The previous example could also be done with queries like these:

```
SELECT listchildren('Top');
listchildren
```

```
-----
Child1
Child2
Child3
(3 rows)
```

```
SELECT name, listchildren(name) FROM nodes;
```

```
name | listchildren
-----+-----
Top   | Child1
Top   | Child2
Top   | Child3
Child1 | SubChild1
Child1 | SubChild2
(5 rows)
```

In the last `SELECT`, notice that no output row appears for `Child2`, `Child3`, etc. This happens because `listchildren` returns an empty set for those arguments, so no result rows are generated. This is the same behavior as we got from an inner join to the function result when using the `LATERAL` syntax.

Postgres Pro's behavior for a set-returning function in a query's select list is almost exactly the same as if the set-returning function had been written in a `LATERAL FROM`-clause item instead. For example,

```
SELECT x, generate_series(1,5) AS g FROM tab;
```

is almost equivalent to

```
SELECT x, g FROM tab, LATERAL generate_series(1,5) AS g;
```

It would be exactly the same, except that in this specific example, the planner could choose to put `g` on the outside of the nested-loop join, since `g` has no actual lateral dependency on `tab`. That would result in a different output row order. Set-returning functions in the select list are always evaluated as though they are on the inside of a nested-loop join with the rest of the `FROM` clause, so that the function(s) are run to completion before the next row from the `FROM` clause is considered.

If there is more than one set-returning function in the query's select list, the behavior is similar to what you get from putting the functions into a single `LATERAL ROWS FROM(...)` `FROM`-clause item. For each row from the underlying query, there is an output row using the first result from each function, then an output row using the second result, and so on. If some of the set-returning functions produce fewer outputs than others, null values are substituted for the missing data, so that the total number of rows emitted for one underlying row is the same as for the set-returning function that produced the most outputs. Thus the set-returning functions run “in lockstep” until they are all exhausted, and then execution continues with the next underlying row.

Set-returning functions can be nested in a select list, although that is not allowed in `FROM`-clause items. In such cases, each level of nesting is treated separately, as though it were a separate `LATERAL ROWS FROM(...)` item. For example, in

```
SELECT srf1(srf2(x), srf3(y)), srf4(srf5(z)) FROM tab;
```

the set-returning functions `srf2`, `srf3`, and `srf5` would be run in lockstep for each row of `tab`, and then `srf1` and `srf4` would be applied in lockstep to each row produced by the lower functions.

Set-returning functions cannot be used within conditional-evaluation constructs, such as `CASE` or `COALESCE`. For example, consider

```
SELECT x, CASE WHEN x > 0 THEN generate_series(1, 5) ELSE 0 END FROM tab;
```

It might seem that this should produce five repetitions of input rows that have `x > 0`, and a single repetition of those that do not; but actually, because `generate_series(1, 5)` would be run in an implicit `LATERAL FROM` item before the `CASE` expression is ever evaluated, it would produce five repetitions of every input row. To reduce confusion, such cases produce a parse-time error instead.

Note

If a function's last command is `INSERT`, `UPDATE`, `DELETE`, or `MERGE` with `RETURNING`, that command will always be executed to completion, even if the function is not declared with `SETOF` or the calling query does not fetch all the result rows. Any extra rows produced by the `RETURNING` clause are silently dropped, but the commanded table modifications still happen (and are all completed before returning from the function).

Note

Before Postgres Pro 10, putting more than one set-returning function in the same select list did not behave very sensibly unless they always produced equal numbers of rows. Otherwise, what you got was a number of output rows equal to the least common multiple of the numbers of rows produced by the set-returning functions. Also, nested set-returning functions did not work as described above; instead, a set-returning function could have at most one set-returning argument, and each nest of set-returning functions was run independently. Also, conditional execution (set-returning functions inside `CASE` etc.) was previously allowed, complicating things even more. Use of the `LATERAL` syntax is recommended when writing queries that need to work in older Postgres Pro versions, because that will give consistent results across different versions. If you have a query that is relying on conditional execution of a set-returning function, you may be able to fix it by moving the conditional test into a custom set-returning function. For example,

```
SELECT x, CASE WHEN y > 0 THEN generate_series(1, z) ELSE 5 END FROM tab;
```


could become

```
CREATE FUNCTION case_generate_series(cond bool, start int, fin int, els int)
  RETURNS SETOF int AS $$
BEGIN
  IF cond THEN
    RETURN QUERY SELECT generate_series(start, fin);
  ELSE
    RETURN QUERY SELECT els;
  END IF;
END$$ LANGUAGE plpgsql;

SELECT x, case_generate_series(y > 0, 1, z, 5) FROM tab;
This formulation will work the same in all versions of Postgres Pro.
```

39.5.10. SQL Functions Returning TABLE

There is another way to declare a function as returning a set, which is to use the syntax `RETURNS TABLE(columns)`. This is equivalent to using one or more `OUT` parameters plus marking the function as returning `SETOF record` (or `SETOF` a single output parameter's type, as appropriate). This notation is specified in recent versions of the SQL standard, and thus may be more portable than using `SETOF`.

For example, the preceding sum-and-product example could also be done this way:

```
CREATE FUNCTION sum_n_product_with_tab (x int)
  RETURNS TABLE(sum int, product int) AS $$
  SELECT $1 + tab.y, $1 * tab.y FROM tab;
$$ LANGUAGE SQL;
```

It is not allowed to use explicit `OUT` or `INOUT` parameters with the `RETURNS TABLE` notation — you must put all the output columns in the `TABLE` list.

39.5.11. Polymorphic SQL Functions

SQL functions can be declared to accept and return the polymorphic types described in [Section 39.2.5](#). Here is a polymorphic function `make_array` that builds up an array from two arbitrary data type elements:

```
CREATE FUNCTION make_array(anelement, anelement) RETURNS anyarray AS $$
  SELECT ARRAY[$1, $2];
$$ LANGUAGE SQL;

SELECT make_array(1, 2) AS intarray, make_array('a'::text, 'b') AS textarray;
 intarray | textarray
-----+-----
 {1,2}   | {a,b}
(1 row)
```

Notice the use of the typecast `'a'::text` to specify that the argument is of type `text`. This is required if the argument is just a string literal, since otherwise it would be treated as type `unknown`, and array of `unknown` is not a valid type. Without the typecast, you will get errors like this:

```
ERROR: could not determine polymorphic type because input has type unknown
```

With `make_array` declared as above, you must provide two arguments that are of exactly the same data type; the system will not attempt to resolve any type differences. Thus for example this does not work:

```
SELECT make_array(1, 2.5) AS numericarray;
ERROR: function make_array(integer, numeric) does not exist
```

An alternative approach is to use the “common” family of polymorphic types, which allows the system to try to identify a suitable common type:

```
CREATE FUNCTION make_array2(anycompatible, anycompatible)
RETURNS anycompatiblearray AS $$
    SELECT ARRAY[$1, $2];
$$ LANGUAGE SQL;
```

```
SELECT make_array2(1, 2.5) AS numericarray;
numericarray
-----
{1,2.5}
(1 row)
```

Because the rules for common type resolution default to choosing type `text` when all inputs are of unknown types, this also works:

```
SELECT make_array2('a', 'b') AS textarray;
textarray
-----
{a,b}
(1 row)
```

It is permitted to have polymorphic arguments with a fixed return type, but the converse is not. For example:

```
CREATE FUNCTION is_greater(anelement, anelement) RETURNS boolean AS $$
    SELECT $1 > $2;
$$ LANGUAGE SQL;
```

```
SELECT is_greater(1, 2);
is_greater
-----
f
(1 row)
```

```
CREATE FUNCTION invalid_func() RETURNS anelement AS $$
    SELECT 1;
$$ LANGUAGE SQL;
ERROR:  cannot determine result data type
DETAIL:  A result of type anelement requires at least one input of type anelement,
anyarray, anynonarray, anyenum, or anyrange.
```

Polymorphism can be used with functions that have output arguments. For example:

```
CREATE FUNCTION dup (f1 anelement, OUT f2 anelement, OUT f3 anyarray)
AS 'select $1, array[$1,$1]' LANGUAGE SQL;
```

```
SELECT * FROM dup(22);
 f2 |   f3
----+-----
 22 | {22,22}
(1 row)
```

Polymorphism can also be used with variadic functions. For example:

```
CREATE FUNCTION anyleast (VARIADIC anyarray) RETURNS anelement AS $$
    SELECT min($1[i]) FROM generate_subscripts($1, 1) g(i);
$$ LANGUAGE SQL;
```

```
SELECT anyleast(10, -1, 5, 4);
anyleast
-----
-1
```

```
(1 row)

SELECT anyleast('abc'::text, 'def');
 anyleast
-----
 abc
(1 row)

CREATE FUNCTION concat_values(text, VARIADIC anyarray) RETURNS text AS $$
    SELECT array_to_string($2, $1);
$$ LANGUAGE SQL;

SELECT concat_values('|', 1, 4, 2);
 concat_values
-----
 1|4|2
(1 row)
```

39.5.12. SQL Functions with Collations

When an SQL function has one or more parameters of collatable data types, a collation is identified for each function call depending on the collations assigned to the actual arguments, as described in [Section 23.2](#). If a collation is successfully identified (i.e., there are no conflicts of implicit collations among the arguments) then all the collatable parameters are treated as having that collation implicitly. This will affect the behavior of collation-sensitive operations within the function. For example, using the `anyleast` function described above, the result of

```
SELECT anyleast('abc'::text, 'ABC');
```

will depend on the database's default collation. In `C` locale the result will be `ABC`, but in many other locales it will be `abc`. The collation to use can be forced by adding a `COLLATE` clause to any of the arguments, for example

```
SELECT anyleast('abc'::text, 'ABC' COLLATE "C");
```

Alternatively, if you wish a function to operate with a particular collation regardless of what it is called with, insert `COLLATE` clauses as needed in the function definition. This version of `anyleast` would always use `en_US` locale to compare strings:

```
CREATE FUNCTION anyleast (VARIADIC anyarray) RETURNS anyelement AS $$
    SELECT min($1[i] COLLATE "en_US") FROM generate_subscripts($1, 1) g(i);
$$ LANGUAGE SQL;
```

But note that this will throw an error if applied to a non-collatable data type.

If no common collation can be identified among the actual arguments, then an SQL function treats its parameters as having their data types' default collation (which is usually the database's default collation, but could be different for parameters of domain types).

The behavior of collatable parameters can be thought of as a limited form of polymorphism, applicable only to textual data types.

39.6. Function Overloading

More than one function can be defined with the same SQL name, so long as the arguments they take are different. In other words, function names can be *overloaded*. Whether or not you use it, this capability entails security precautions when calling functions in databases where some users mistrust other users; see [Section 10.3](#). When a query is executed, the server will determine which function to call from the data types and the number of the provided arguments. Overloading can also be used to simulate functions with a variable number of arguments, up to a finite maximum number.

When creating a family of overloaded functions, one should be careful not to create ambiguities. For instance, given the functions:

```
CREATE FUNCTION test(int, real) RETURNS ...
CREATE FUNCTION test(smallint, double precision) RETURNS ...
```

it is not immediately clear which function would be called with some trivial input like `test(1, 1.5)`. The currently implemented resolution rules are described in [Chapter 10](#), but it is unwise to design a system that subtly relies on this behavior.

A function that takes a single argument of a composite type should generally not have the same name as any attribute (field) of that type. Recall that `attribute(table)` is considered equivalent to `table.attribute`. In the case that there is an ambiguity between a function on a composite type and an attribute of the composite type, the attribute will always be used. It is possible to override that choice by schema-qualifying the function name (that is, `schema.func(table)`) but it's better to avoid the problem by not choosing conflicting names.

Another possible conflict is between variadic and non-variadic functions. For instance, it is possible to create both `foo(numeric)` and `foo(VARIADIC numeric[])`. In this case it is unclear which one should be matched to a call providing a single numeric argument, such as `foo(10.1)`. The rule is that the function appearing earlier in the search path is used, or if the two functions are in the same schema, the non-variadic one is preferred.

When overloading C-language functions, there is an additional constraint: The C name of each function in the family of overloaded functions must be different from the C names of all other functions, either internal or dynamically loaded. If this rule is violated, the behavior is not portable. You might get a runtime linker error, or one of the functions will get called (usually the internal one). The alternative form of the `AS` clause for the SQL `CREATE FUNCTION` command decouples the SQL function name from the function name in the C source code. For instance:

```
CREATE FUNCTION test(int) RETURNS int
    AS 'filename', 'test_1arg'
    LANGUAGE C;
CREATE FUNCTION test(int, int) RETURNS int
    AS 'filename', 'test_2arg'
    LANGUAGE C;
```

The names of the C functions here reflect one of many possible conventions.

39.7. Function Volatility Categories

Every function has a *volatility* classification, with the possibilities being `VOLATILE`, `STABLE`, or `IMMUTABLE`. `VOLATILE` is the default if the `CREATE FUNCTION` command does not specify a category. The volatility category is a promise to the optimizer about the behavior of the function:

- A `VOLATILE` function can do anything, including modifying the database. It can return different results on successive calls with the same arguments. The optimizer makes no assumptions about the behavior of such functions. A query using a volatile function will re-evaluate the function at every row where its value is needed.
- A `STABLE` function cannot modify the database and is guaranteed to return the same results given the same arguments for all rows within a single statement. This category allows the optimizer to optimize multiple calls of the function to a single call. In particular, it is safe to use an expression containing such a function in an index scan condition. (Since an index scan will evaluate the comparison value only once, not once at each row, it is not valid to use a `VOLATILE` function in an index scan condition.)
- An `IMMUTABLE` function cannot modify the database and is guaranteed to return the same results given the same arguments forever. This category allows the optimizer to pre-evaluate the function when a query calls it with constant arguments. For example, a query like `SELECT ... WHERE x = 2 + 2` can be simplified on sight to `SELECT ... WHERE x = 4`, because the function underlying the integer addition operator is marked `IMMUTABLE`.

For best optimization results, you should label your functions with the strictest volatility category that is valid for them.

Any function with side-effects *must* be labeled `VOLATILE`, so that calls to it cannot be optimized away. Even a function with no side-effects needs to be labeled `VOLATILE` if its value can change within a single query; some examples are `random()`, `currval()`, `timeofday()`.

Another important example is that the `current_timestamp` family of functions qualify as `STABLE`, since their values do not change within a transaction.

There is relatively little difference between `STABLE` and `IMMUTABLE` categories when considering simple interactive queries that are planned and immediately executed: it doesn't matter a lot whether a function is executed once during planning or once during query execution startup. But there is a big difference if the plan is saved and reused later. Labeling a function `IMMUTABLE` when it really isn't might allow it to be prematurely folded to a constant during planning, resulting in a stale value being re-used during subsequent uses of the plan. This is a hazard when using prepared statements or when using function languages that cache plans (such as PL/pgSQL).

For functions written in SQL or in any of the standard procedural languages, there is a second important property determined by the volatility category, namely the visibility of any data changes that have been made by the SQL command that is calling the function. A `VOLATILE` function will see such changes, a `STABLE` or `IMMUTABLE` function will not. This behavior is implemented using the snapshotting behavior of MVCC (see [Chapter 13](#)): `STABLE` and `IMMUTABLE` functions use a snapshot established as of the start of the calling query, whereas `VOLATILE` functions obtain a fresh snapshot at the start of each query they execute.

Note

Functions written in C can manage snapshots however they want, but it's usually a good idea to make C functions work this way too.

Because of this snapshotting behavior, a function containing only `SELECT` commands can safely be marked `STABLE`, even if it selects from tables that might be undergoing modifications by concurrent queries. Postgres Pro will execute all commands of a `STABLE` function using the snapshot established for the calling query, and so it will see a fixed view of the database throughout that query.

The same snapshotting behavior is used for `SELECT` commands within `IMMUTABLE` functions. It is generally unwise to select from database tables within an `IMMUTABLE` function at all, since the immutability will be broken if the table contents ever change. However, Postgres Pro does not enforce that you do not do that.

A common error is to label a function `IMMUTABLE` when its results depend on a configuration parameter. For example, a function that manipulates timestamps might well have results that depend on the [Time-Zone](#) setting. For safety, such functions should be labeled `STABLE` instead.

Note

Postgres Pro requires that `STABLE` and `IMMUTABLE` functions contain no SQL commands other than `SELECT` to prevent data modification. (This is not a completely bulletproof test, since such functions could still call `VOLATILE` functions that modify the database. If you do that, you will find that the `STABLE` or `IMMUTABLE` function does not notice the database changes applied by the called function, since they are hidden from its snapshot.)

39.8. Procedural Language Functions

Postgres Pro allows user-defined functions to be written in other languages besides SQL and C. These other languages are generically called *procedural languages* (PLs). Procedural languages aren't built into the Postgres Pro server; they are offered by loadable modules. See [Chapter 43](#) and following chapters for more information.

39.9. Internal Functions

Internal functions are functions written in C that have been statically linked into the Postgres Pro server. The “body” of the function definition specifies the C-language name of the function, which need not be the same as the name being declared for SQL use. (For reasons of backward compatibility, an empty body is accepted as meaning that the C-language function name is the same as the SQL name.)

Normally, all internal functions present in the server are declared during the initialization of the database cluster (see [Section 18.2](#)), but a user could use `CREATE FUNCTION` to create additional alias names for an internal function. Internal functions are declared in `CREATE FUNCTION` with language name `internal`. For instance, to create an alias for the `sqrt` function:

```
CREATE FUNCTION square_root(double precision) RETURNS double precision
AS 'dsqrt'
LANGUAGE internal
STRICT;
```

(Most internal functions expect to be declared “strict”.)

Note

Not all “predefined” functions are “internal” in the above sense. Some predefined functions are written in SQL.

39.10. C-Language Functions

User-defined functions can be written in C (or a language that can be made compatible with C, such as C++). Such functions are compiled into dynamically loadable objects (also called shared libraries) and are loaded by the server on demand. The dynamic loading feature is what distinguishes “C language” functions from “internal” functions — the actual coding conventions are essentially the same for both. (Hence, the standard internal function library is a rich source of coding examples for user-defined C functions.)

Currently only one calling convention is used for C functions (“version 1”). Support for that calling convention is indicated by writing a `PG_FUNCTION_INFO_V1()` macro call for the function, as illustrated below.

39.10.1. Dynamic Loading

The first time a user-defined function in a particular loadable object file is called in a session, the dynamic loader loads that object file into memory so that the function can be called. The `CREATE FUNCTION` for a user-defined C function must therefore specify two pieces of information for the function: the name of the loadable object file, and the C name (link symbol) of the specific function to call within that object file. If the C name is not explicitly specified then it is assumed to be the same as the SQL function name.

The following algorithm is used to locate the shared object file based on the name given in the `CREATE FUNCTION` command:

1. If the name is an absolute path, the given file is loaded.
2. If the name starts with the string `$libdir`, that part is replaced by the Postgres Pro package library directory name, which is determined at build time.
3. If the name does not contain a directory part, the file is searched for in the path specified by the configuration variable `dynamic_library_path`.
4. Otherwise (the file was not found in the path, or it contains a non-absolute directory part), the dynamic loader will try to take the name as given, which will most likely fail. (It is unreliable to depend on the current working directory.)

If this sequence does not work, the platform-specific shared library file name extension (often `.so`) is appended to the given name and this sequence is tried again. If that fails as well, the load will fail.

It is recommended to locate shared libraries either relative to `$libdir` or through the dynamic library path. This simplifies version upgrades if the new installation is at a different location. The actual directory that `$libdir` stands for can be found out with the command `pg_config --pkglibdir`.

The user ID the Postgres Pro server runs as must be able to traverse the path to the file you intend to load. Making the file or a higher-level directory not readable and/or not executable by the postgres user is a common mistake.

In any case, the file name that is given in the `CREATE FUNCTION` command is recorded literally in the system catalogs, so if the file needs to be loaded again the same procedure is applied.

Note

Postgres Pro will not compile a C function automatically. The object file must be compiled before it is referenced in a `CREATE FUNCTION` command. See [Section 39.10.5](#) for additional information.

To ensure that a dynamically loaded object file is not loaded into an incompatible server, Postgres Pro checks that the file contains a “magic block” with the appropriate contents. This allows the server to detect obvious incompatibilities, such as code compiled for a different major version of Postgres Pro. To include a magic block, write this in one (and only one) of the module source files, after having included the header `fmgr.h`:

```
PG_MODULE_MAGIC;
```

After it is used for the first time, a dynamically loaded object file is retained in memory. Future calls in the same session to the function(s) in that file will only incur the small overhead of a symbol table lookup. If you need to force a reload of an object file, for example after recompiling it, begin a fresh session.

Optionally, a dynamically loaded file can contain an initialization function. If the file includes a function named `_PG_init`, that function will be called immediately after loading the file. The function receives no parameters and should return void. There is presently no way to unload a dynamically loaded file.

39.10.2. Base Types in C-Language Functions

To know how to write C-language functions, you need to know how Postgres Pro internally represents base data types and how they can be passed to and from functions. Internally, Postgres Pro regards a base type as a “blob of memory”. The user-defined functions that you define over a type in turn define the way that Postgres Pro can operate on it. That is, Postgres Pro will only store and retrieve the data from disk and use your user-defined functions to input, process, and output the data.

Base types can have one of three internal formats:

- pass by value, fixed-length
- pass by reference, fixed-length
- pass by reference, variable-length

By-value types can only be 1, 2, or 4 bytes in length (also 8 bytes, if `sizeof(Datum)` is 8 on your machine). You should be careful to define your types such that they will be the same size (in bytes) on all architectures. For example, the `long` type is dangerous because it is 4 bytes on some machines and 8 bytes on others, whereas `int` type is 4 bytes on most Unix machines. A reasonable implementation of the `int4` type on Unix machines might be:

```
/* 4-byte integer, passed by value */
typedef int int4;
```

(The actual Postgres Pro C code calls this type `int32`, because it is a convention in C that `intXX` means `XX bits`. Note therefore also that the C type `int8` is 1 byte in size. The SQL type `int8` is called `int64` in C. See also [Table 39.2](#).)

On the other hand, fixed-length types of any size can be passed by-reference. For example, here is a sample implementation of a Postgres Pro type:

```
/* 16-byte structure, passed by reference */
typedef struct
{
    double   x, y;
} Point;
```

Only pointers to such types can be used when passing them in and out of Postgres Pro functions. To return a value of such a type, allocate the right amount of memory with `palloc`, fill in the allocated memory, and return a pointer to it. (Also, if you just want to return the same value as one of your input arguments that's of the same data type, you can skip the extra `palloc` and just return the pointer to the input value.)

Finally, all variable-length types must also be passed by reference. All variable-length types must begin with an opaque length field of exactly 4 bytes, which will be set by `SET_VARSIZE`; never set this field directly! All data to be stored within that type must be located in the memory immediately following that length field. The length field contains the total length of the structure, that is, it includes the size of the length field itself.

Another important point is to avoid leaving any uninitialized bits within data type values; for example, take care to zero out any alignment padding bytes that might be present in structs. Without this, logically-equivalent constants of your data type might be seen as unequal by the planner, leading to inefficient (though not incorrect) plans.

Warning

Never modify the contents of a pass-by-reference input value. If you do so you are likely to corrupt on-disk data, since the pointer you are given might point directly into a disk buffer. The sole exception to this rule is explained in [Section 39.12](#).

As an example, we can define the type `text` as follows:

```
typedef struct {
    int32 length;
    char data[FLEXIBLE_ARRAY_MEMBER];
} text;
```

The `[FLEXIBLE_ARRAY_MEMBER]` notation means that the actual length of the data part is not specified by this declaration.

When manipulating variable-length types, we must be careful to allocate the correct amount of memory and set the length field correctly. For example, if we wanted to store 40 bytes in a `text` structure, we might use a code fragment like this:

```
#include "postgres.h"
...
char buffer[40]; /* our source data */
...
text *destination = (text *) palloc(VARHDRSZ + 40);
SET_VARSIZE(destination, VARHDRSZ + 40);
memcpy(destination->data, buffer, 40);
...
```

`VARHDRSZ` is the same as `sizeof(int32)`, but it's considered good style to use the macro `VARHDRSZ` to refer to the size of the overhead for a variable-length type. Also, the length field *must* be set using the `SET_VARSIZE` macro, not by simple assignment.

[Table 39.2](#) shows the C types corresponding to many of the built-in SQL data types of Postgres Pro. The “Defined In” column gives the header file that needs to be included to get the type definition. (The actual

definition might be in a different file that is included by the listed file. It is recommended that users stick to the defined interface.) Note that you should always include `postgres.h` first in any source file of server code, because it declares a number of things that you will need anyway, and because including other headers first can cause portability issues.

Table 39.2. Equivalent C Types for Built-in SQL Types

SQL Type	C Type	Defined In
boolean	bool	postgres.h (maybe compiler built-in)
box	BOX*	utils/geo_decls.h
bytea	bytea*	postgres.h
"char"	char	(compiler built-in)
character	BpChar*	postgres.h
cid	CommandId	postgres.h
date	DateADT	utils/date.h
float4 (real)	float4	postgres.h
float8 (double precision)	float8	postgres.h
int2 (smallint)	int16	postgres.h
int4 (integer)	int32	postgres.h
int8 (bigint)	int64	postgres.h
interval	Interval*	datatype/timestamp.h
lseg	LSEG*	utils/geo_decls.h
name	Name	postgres.h
numeric	Numeric	utils/numeric.h
oid	Oid	postgres.h
oidvector	oidvector*	postgres.h
path	PATH*	utils/geo_decls.h
point	POINT*	utils/geo_decls.h
regproc	RegProcedure	postgres.h
text	text*	postgres.h
tid	ItemPointer	storage/itemptr.h
time	TimeADT	utils/date.h
time with time zone	TimeTzADT	utils/date.h
timestamp	Timestamp	datatype/timestamp.h
timestamp with time zone	TimestampTz	datatype/timestamp.h
varchar	VarChar*	postgres.h
xid	TransactionId	postgres.h

Now that we've gone over all of the possible structures for base types, we can show some examples of real functions.

39.10.3. Version 1 Calling Conventions

The version-1 calling convention relies on macros to suppress most of the complexity of passing arguments and results. The C declaration of a version-1 function is always:

```
Datum funcname(PG_FUNCTION_ARGS)
```

In addition, the macro call:

```
PG_FUNCTION_INFO_V1(funcname);
```

must appear in the same source file. (Conventionally, it's written just before the function itself.) This macro call is not needed for internal-language functions, since Postgres Pro assumes that all internal functions use the version-1 convention. It is, however, required for dynamically-loaded functions.

In a version-1 function, each actual argument is fetched using a `PG_GETARG_xxx()` macro that corresponds to the argument's data type. (In non-strict functions there needs to be a previous check about argument null-ness using `PG_ARGISNULL()`; see below.) The result is returned using a `PG_RETURN_xxx()` macro for the return type. `PG_GETARG_xxx()` takes as its argument the number of the function argument to fetch, where the count starts at 0. `PG_RETURN_xxx()` takes as its argument the actual value to return.

Here are some examples using the version-1 calling convention:

```
#include "postgres.h"
#include <string.h>
#include "fmgr.h"
#include "utils/geo_decls.h"
#include "varatt.h"

PG_MODULE_MAGIC;

/* by value */

PG_FUNCTION_INFO_V1(add_one);

Datum
add_one(PG_FUNCTION_ARGS)
{
    int32    arg = PG_GETARG_INT32(0);

    PG_RETURN_INT32(arg + 1);
}

/* by reference, fixed length */

PG_FUNCTION_INFO_V1(add_one_float8);

Datum
add_one_float8(PG_FUNCTION_ARGS)
{
    /* The macros for FLOAT8 hide its pass-by-reference nature. */
    float8   arg = PG_GETARG_FLOAT8(0);

    PG_RETURN_FLOAT8(arg + 1.0);
}

PG_FUNCTION_INFO_V1(makepoint);

Datum
makepoint(PG_FUNCTION_ARGS)
{
    /* Here, the pass-by-reference nature of Point is not hidden. */
    Point    *pointx = PG_GETARG_POINT_P(0);
    Point    *pointy = PG_GETARG_POINT_P(1);
    Point    *new_point = (Point *) palloc(sizeof(Point));
```

```

    new_point->x = pointx->x;
    new_point->y = pointy->y;

    PG_RETURN_POINT_P(new_point);
}

/* by reference, variable length */

PG_FUNCTION_INFO_V1(copytext);

Datum
copytext(PG_FUNCTION_ARGS)
{
    text      *t = PG_GETARG_TEXT_PP(0);

    /*
     * VARSIZE_ANY_EXHDR is the size of the struct in bytes, minus the
     * VARHDRSZ or VARHDRSZ_SHORT of its header. Construct the copy with a
     * full-length header.
     */
    text      *new_t = (text *) palloc(VARSIZE_ANY_EXHDR(t) + VARHDRSZ);
    SET_VARSIZE(new_t, VARSIZE_ANY_EXHDR(t) + VARHDRSZ);

    /*
     * VARDATA is a pointer to the data region of the new struct. The source
     * could be a short datum, so retrieve its data through VARDATA_ANY.
     */
    memcpy(VARDATA(new_t),          /* destination */
           VARDATA_ANY(t),          /* source */
           VARSIZE_ANY_EXHDR(t));   /* how many bytes */
    PG_RETURN_TEXT_P(new_t);
}

PG_FUNCTION_INFO_V1(concat_text);

Datum
concat_text(PG_FUNCTION_ARGS)
{
    text *arg1 = PG_GETARG_TEXT_PP(0);
    text *arg2 = PG_GETARG_TEXT_PP(1);
    int32 arg1_size = VARSIZE_ANY_EXHDR(arg1);
    int32 arg2_size = VARSIZE_ANY_EXHDR(arg2);
    int32 new_text_size = arg1_size + arg2_size + VARHDRSZ;
    text *new_text = (text *) palloc(new_text_size);

    SET_VARSIZE(new_text, new_text_size);
    memcpy(VARDATA(new_text), VARDATA_ANY(arg1), arg1_size);
    memcpy(VARDATA(new_text) + arg1_size, VARDATA_ANY(arg2), arg2_size);
    PG_RETURN_TEXT_P(new_text);
}

```

Supposing that the above code has been prepared in file `funcs.c` and compiled into a shared object, we could define the functions to Postgres Pro with commands like this:

```

CREATE FUNCTION add_one(integer) RETURNS integer
AS 'DIRECTORY/funcs', 'add_one'
LANGUAGE C STRICT;

```

```
-- note overloading of SQL function name "add_one"
CREATE FUNCTION add_one(double precision) RETURNS double precision
    AS 'DIRECTORY/funcs', 'add_one_float8'
    LANGUAGE C STRICT;

CREATE FUNCTION makepoint(point, point) RETURNS point
    AS 'DIRECTORY/funcs', 'makepoint'
    LANGUAGE C STRICT;

CREATE FUNCTION copytext(text) RETURNS text
    AS 'DIRECTORY/funcs', 'copytext'
    LANGUAGE C STRICT;

CREATE FUNCTION concat_text(text, text) RETURNS text
    AS 'DIRECTORY/funcs', 'concat_text'
    LANGUAGE C STRICT;
```

Here, *DIRECTORY* stands for the directory of the shared library file (for instance the Postgres Pro tutorial directory, which contains the code for the examples used in this section). (Better style would be to use just 'funcs' in the AS clause, after having added *DIRECTORY* to the search path. In any case, we can omit the system-specific extension for a shared library, commonly .so.)

Notice that we have specified the functions as “strict”, meaning that the system should automatically assume a null result if any input value is null. By doing this, we avoid having to check for null inputs in the function code. Without this, we'd have to check for null values explicitly, using `PG_ARGISNULL()`.

The macro `PG_ARGISNULL(n)` allows a function to test whether each input is null. (Of course, doing this is only necessary in functions not declared “strict”.) As with the `PG_GETARG_xxx()` macros, the input arguments are counted beginning at zero. Note that one should refrain from executing `PG_GETARG_xxx()` until one has verified that the argument isn't null. To return a null result, execute `PG_RETURN_NULL()`; this works in both strict and nonstrict functions.

At first glance, the version-1 coding conventions might appear to be just pointless obscurantism, compared to using plain C calling conventions. They do however allow us to deal with NULLable arguments/return values, and “toasted” (compressed or out-of-line) values.

Other options provided by the version-1 interface are two variants of the `PG_GETARG_xxx()` macros. The first of these, `PG_GETARG_xxx_COPY()`, guarantees to return a copy of the specified argument that is safe for writing into. (The normal macros will sometimes return a pointer to a value that is physically stored in a table, which must not be written to. Using the `PG_GETARG_xxx_COPY()` macros guarantees a writable result.) The second variant consists of the `PG_GETARG_xxx_SLICE()` macros which take three arguments. The first is the number of the function argument (as above). The second and third are the offset and length of the segment to be returned. Offsets are counted from zero, and a negative length requests that the remainder of the value be returned. These macros provide more efficient access to parts of large values in the case where they have storage type “external”. (The storage type of a column can be specified using `ALTER TABLE tablename ALTER COLUMN colname SET STORAGE storagetype`. *storagetype* is one of plain, external, extended, or main.)

Finally, the version-1 function call conventions make it possible to return set results ([Section 39.10.8](#)) and implement trigger functions ([Chapter 40](#)) and procedural-language call handlers ([Chapter 57](#)).

39.10.4. Writing Code

Before we turn to the more advanced topics, we should discuss some coding rules for Postgres Pro C-language functions. While it might be possible to load functions written in languages other than C into Postgres Pro, this is usually difficult (when it is possible at all) because other languages, such as C++, FORTRAN, or Pascal often do not follow the same calling convention as C. That is, other languages do not pass argument and return values between functions in the same way. For this reason, we will assume that your C-language functions are actually written in C.

The basic rules for writing and building C functions are as follows:

- Use `pg_config --includedir-server` to find out where the Postgres Pro server header files are installed on your system (or the system that your users will be running on).
- Compiling and linking your code so that it can be dynamically loaded into Postgres Pro always requires special flags. See [Section 39.10.5](#) for a detailed explanation of how to do it for your particular operating system.
- Remember to define a “magic block” for your shared library, as described in [Section 39.10.1](#).
- When allocating memory, use the Postgres Pro functions `palloc` and `pfree` instead of the corresponding C library functions `malloc` and `free`. The memory allocated by `palloc` will be freed automatically at the end of each transaction, preventing memory leaks.
- Always zero the bytes of your structures using `memset` (or allocate them with `palloc0` in the first place). Even if you assign to each field of your structure, there might be alignment padding (holes in the structure) that contain garbage values. Without this, it's difficult to support hash indexes or hash joins, as you must pick out only the significant bits of your data structure to compute a hash. The planner also sometimes relies on comparing constants via bitwise equality, so you can get undesirable planning results if logically-equivalent values aren't bitwise equal.
- Most of the internal Postgres Pro types are declared in `postgres.h`, while the function manager interfaces (`PG_FUNCTION_ARGS`, etc.) are in `fmgr.h`, so you will need to include at least these two files. For portability reasons it's best to include `postgres.h` *first*, before any other system or user header files. Including `postgres.h` will also include `elog.h` and `palloc.h` for you.
- Symbol names defined within object files must not conflict with each other or with symbols defined in the Postgres Pro server executable. You will have to rename your functions or variables if you get error messages to this effect.

39.10.5. Compiling and Linking Dynamically-Loaded Functions

Before you are able to use your Postgres Pro extension functions written in C, they must be compiled and linked in a special way to produce a file that can be dynamically loaded by the server. To be precise, a *shared library* needs to be created.

For information beyond what is contained in this section you should read the documentation of your operating system, in particular the manual pages for the C compiler, `cc`, and the link editor, `ld`.

Creating shared libraries is generally analogous to linking executables: first the source files are compiled into object files, then the object files are linked together. The object files need to be created as *position-independent code* (PIC), which conceptually means that they can be placed at an arbitrary location in memory when they are loaded by the executable. (Object files intended for executables are usually not compiled that way.) The command to link a shared library contains special flags to distinguish it from linking an executable (at least in theory — on some systems the practice is much uglier).

In the following examples we assume that your source code is in a file `foo.c` and we will create a shared library `foo.so`. The intermediate object file will be called `foo.o` unless otherwise noted. A shared library can contain more than one object file, but we only use one here.

FreeBSD

The compiler flag to create PIC is `-fPIC`. To create shared libraries the compiler flag is `-shared`.

```
cc -fPIC -c foo.c
cc -shared -o foo.so foo.o
```

This is applicable as of version 13.0 of FreeBSD, older versions used the `gcc` compiler.

Linux

The compiler flag to create PIC is `-fPIC`. The compiler flag to create a shared library is `-shared`. A complete example looks like this:

```
cc -fPIC -c foo.c
cc -shared -o foo.so foo.o
```

macOS

Here is an example. It assumes the developer tools are installed.

```
cc -c foo.c
cc -bundle -flat_namespace -undefined suppress -o foo.so foo.o
```

NetBSD

The compiler flag to create PIC is `-fPIC`. For ELF systems, the compiler with the flag `-shared` is used to link shared libraries. On the older non-ELF systems, `ld -Bshareable` is used.

```
gcc -fPIC -c foo.c
gcc -shared -o foo.so foo.o
```

OpenBSD

The compiler flag to create PIC is `-fPIC`. `ld -Bshareable` is used to link shared libraries.

```
gcc -fPIC -c foo.c
ld -Bshareable -o foo.so foo.o
```

Solaris

The compiler flag to create PIC is `-KPIC` with the Sun compiler and `-fPIC` with GCC. To link shared libraries, the compiler option is `-G` with either compiler or alternatively `-shared` with GCC.

```
cc -KPIC -c foo.c
cc -G -o foo.so foo.o
```

or

```
gcc -fPIC -c foo.c
gcc -G -o foo.so foo.o
```

Tip

If this is too complicated for you, you should consider using [GNU Libtool](#), which hides the platform differences behind a uniform interface.

The resulting shared library file can then be loaded into Postgres Pro. When specifying the file name to the `CREATE FUNCTION` command, one must give it the name of the shared library file, not the intermediate object file. Note that the system's standard shared-library extension (usually `.so` or `.sl`) can be omitted from the `CREATE FUNCTION` command, and normally should be omitted for best portability.

Refer back to [Section 39.10.1](#) about where the server expects to find the shared library files.

39.10.6. Composite-Type Arguments

Composite types do not have a fixed layout like C structures. Instances of a composite type can contain null fields. In addition, composite types that are part of an inheritance hierarchy can have different fields than other members of the same inheritance hierarchy. Therefore, Postgres Pro provides a function interface for accessing fields of composite types from C.

Suppose we want to write a function to answer the query:

```
SELECT name, c_overpaid(emp, 1500) AS overpaid
FROM emp
WHERE name = 'Bill' OR name = 'Sam';
```

Using the version-1 calling conventions, we can define `c_overpaid` as:

```
#include "postgres.h"
#include "executor/executor.h" /* for GetAttributeByName() */

PG_MODULE_MAGIC;

PG_FUNCTION_INFO_V1(c_overpaid);

Datum
c_overpaid(PG_FUNCTION_ARGS)
{
    HeapTupleHeader t = PG_GETARG_HEAPTUPLEHEADER(0);
    int32          limit = PG_GETARG_INT32(1);
    bool isnull;
    Datum salary;

    salary = GetAttributeByName(t, "salary", &isnull);
    if (isnull)
        PG_RETURN_BOOL(false);
    /* Alternatively, we might prefer to do PG_RETURN_NULL() for null salary. */

    PG_RETURN_BOOL(DatumGetInt32(salary) > limit);
}
```

`GetAttributeByName` is the Postgres Pro system function that returns attributes out of the specified row. It has three arguments: the argument of type `HeapTupleHeader` passed into the function, the name of the desired attribute, and a return parameter that tells whether the attribute is null. `GetAttributeByName` returns a `Datum` value that you can convert to the proper data type by using the appropriate `DatumGetXXX()` function. Note that the return value is meaningless if the null flag is set; always check the null flag before trying to do anything with the result.

There is also `GetAttributeByNum`, which selects the target attribute by column number instead of name.

The following command declares the function `c_overpaid` in SQL:

```
CREATE FUNCTION c_overpaid(emp, integer) RETURNS boolean
AS 'DIRECTORY/funcs', 'c_overpaid'
LANGUAGE C STRICT;
```

Notice we have used `STRICT` so that we did not have to check whether the input arguments were `NULL`.

39.10.7. Returning Rows (Composite Types)

To return a row or composite-type value from a C-language function, you can use a special API that provides macros and functions to hide most of the complexity of building composite data types. To use this API, the source file must include:

```
#include "funcapi.h"
```

There are two ways you can build a composite data value (henceforth a “tuple”): you can build it from an array of `Datum` values, or from an array of C strings that can be passed to the input conversion functions of the tuple's column data types. In either case, you first need to obtain or construct a `TupleDesc` descriptor for the tuple structure. When working with `Datums`, you pass the `TupleDesc` to `BlessTupleDesc`, and then call `heap_form_tuple` for each row. When working with C strings, you pass the `TupleDesc` to `TupleDescGetAttInMetadata`, and then call `BuildTupleFromCStrings` for each row. In the case of a function returning a set of tuples, the setup steps can all be done once during the first call of the function.

Several helper functions are available for setting up the needed `TupleDesc`. The recommended way to do this in most functions returning composite values is to call:

```
TypeFuncClass get_call_result_type(FunctionCallInfo fcinfo,
                                   Oid *resultTypeId,
                                   TupleDesc *resultTupleDesc)
```

passing the same `fcinfo` struct passed to the calling function itself. (This of course requires that you use the version-1 calling conventions.) `resultTypeId` can be specified as `NULL` or as the address of a local variable to receive the function's result type OID. `resultTupleDesc` should be the address of a local `TupleDesc` variable. Check that the result is `TYPEFUNC_COMPOSITE`; if so, `resultTupleDesc` has been filled with the needed `TupleDesc`. (If it is not, you can report an error along the lines of “function returning record called in context that cannot accept type record”.)

Tip

`get_call_result_type` can resolve the actual type of a polymorphic function result; so it is useful in functions that return scalar polymorphic results, not only functions that return composites. The `resultTypeId` output is primarily useful for functions returning polymorphic scalars.

Note

`get_call_result_type` has a sibling `get_expr_result_type`, which can be used to resolve the expected output type for a function call represented by an expression tree. This can be used when trying to determine the result type from outside the function itself. There is also `get_func_result_type`, which can be used when only the function's OID is available. However these functions are not able to deal with functions declared to return `record`, and `get_func_result_type` cannot resolve polymorphic types, so you should preferentially use `get_call_result_type`.

Older, now-deprecated functions for obtaining `TupleDesc`s are:

```
TupleDesc RelationNameGetTupleDesc(const char *relname)
```

to get a `TupleDesc` for the row type of a named relation, and:

```
TupleDesc TypeGetTupleDesc(Oid typeoid, List *colaliases)
```

to get a `TupleDesc` based on a type OID. This can be used to get a `TupleDesc` for a base or composite type. It will not work for a function that returns `record`, however, and it cannot resolve polymorphic types.

Once you have a `TupleDesc`, call:

```
TupleDesc BlessTupleDesc(TupleDesc tupdesc)
```

if you plan to work with Datums, or:

```
AttInMetadata *TupleDescGetAttInMetadata(TupleDesc tupdesc)
```

if you plan to work with C strings. If you are writing a function returning set, you can save the results of these functions in the `FuncCallContext` structure — use the `tuple_desc` or `attinmeta` field respectively.

When working with Datums, use:

```
HeapTuple heap_form_tuple(TupleDesc tupdesc, Datum *values, bool *isnull)
```

to build a `HeapTuple` given user data in Datum form.

When working with C strings, use:

```
HeapTuple BuildTupleFromCStrings(AttInMetadata *attinmeta, char **values)
```

to build a `HeapTuple` given user data in C string form. `values` is an array of C strings, one for each attribute of the return row. Each C string should be in the form expected by the input function of the

attribute data type. In order to return a null value for one of the attributes, the corresponding pointer in the *values* array should be set to `NULL`. This function will need to be called again for each row you return.

Once you have built a tuple to return from your function, it must be converted into a `Datum`. Use:

```
HeapTupleGetDatum(HeapTuple tuple)
```

to convert a `HeapTuple` into a valid `Datum`. This `Datum` can be returned directly if you intend to return just a single row, or it can be used as the current return value in a set-returning function.

An example appears in the next section.

39.10.8. Returning Sets

C-language functions have two options for returning sets (multiple rows). In one method, called *ValuePerCall* mode, a set-returning function is called repeatedly (passing the same arguments each time) and it returns one new row on each call, until it has no more rows to return and signals that by returning `NULL`. The set-returning function (SRF) must therefore save enough state across calls to remember what it was doing and return the correct next item on each call. In the other method, called *Materialize* mode, an SRF fills and returns a tuplestore object containing its entire result; then only one call occurs for the whole result, and no inter-call state is needed.

When using *ValuePerCall* mode, it is important to remember that the query is not guaranteed to be run to completion; that is, due to options such as `LIMIT`, the executor might stop making calls to the set-returning function before all rows have been fetched. This means it is not safe to perform cleanup activities in the last call, because that might not ever happen. It's recommended to use *Materialize* mode for functions that need access to external resources, such as file descriptors.

The remainder of this section documents a set of helper macros that are commonly used (though not required to be used) for SRFs using *ValuePerCall* mode.

To use the *ValuePerCall* support macros described here, include `funcapi.h`. These macros work with a structure `FuncCallContext` that contains the state that needs to be saved across calls. Within the calling SRF, `fcinfo->flinfo->fn_extra` is used to hold a pointer to `FuncCallContext` across calls. The macros automatically fill that field on first use, and expect to find the same pointer there on subsequent uses.

```
typedef struct FuncCallContext
{
    /*
     * Number of times we've been called before
     */
    /* call_cntr is initialized to 0 for you by SRF_FIRSTCALL_INIT(), and
     * incremented for you every time SRF_RETURN_NEXT() is called.
     */
    uint64 call_cntr;

    /*
     * OPTIONAL maximum number of calls
     */
    /* max_calls is here for convenience only and setting it is optional.
     * If not set, you must provide alternative means to know when the
     * function is done.
     */
    uint64 max_calls;

    /*
     * OPTIONAL pointer to miscellaneous user-provided context information
     */
    /* user_fctx is for use as a pointer to your own data to retain
     * arbitrary context information between calls of your function.
     */
}
```

```

void *user_fctx;

/*
 * OPTIONAL pointer to struct containing attribute type input metadata
 *
 * attinmeta is for use when returning tuples (i.e., composite data types)
 * and is not used when returning base data types. It is only needed
 * if you intend to use BuildTupleFromCStrings() to create the return
 * tuple.
 */
AttInMetadata *attinmeta;

/*
 * memory context used for structures that must live for multiple calls
 *
 * multi_call_memory_ctx is set by SRF_FIRSTCALL_INIT() for you, and used
 * by SRF_RETURN_DONE() for cleanup. It is the most appropriate memory
 * context for any memory that is to be reused across multiple calls
 * of the SRF.
 */
MemoryContext multi_call_memory_ctx;

/*
 * OPTIONAL pointer to struct containing tuple description
 *
 * tuple_desc is for use when returning tuples (i.e., composite data types)
 * and is only needed if you are going to build the tuples with
 * heap_form_tuple() rather than with BuildTupleFromCStrings(). Note that
 * the TupleDesc pointer stored here should usually have been run through
 * BlessTupleDesc() first.
 */
TupleDesc tuple_desc;

} FuncCallContext;

```

The macros to be used by an SRF using this infrastructure are:

`SRF_IS_FIRSTCALL()`

Use this to determine if your function is being called for the first or a subsequent time. On the first call (only), call:

`SRF_FIRSTCALL_INIT()`

to initialize the `FuncCallContext`. On every function call, including the first, call:

`SRF_PERCALL_SETUP()`

to set up for using the `FuncCallContext`.

If your function has data to return in the current call, use:

`SRF_RETURN_NEXT(funcctx, result)`

to return it to the caller. (`result` must be of type `Datum`, either a single value or a tuple prepared as described above.) Finally, when your function is finished returning data, use:

`SRF_RETURN_DONE(funcctx)`

to clean up and end the SRF.

The memory context that is current when the SRF is called is a transient context that will be cleared between calls. This means that you do not need to call `pfree` on everything you allocated using `pal-`

loc; it will go away anyway. However, if you want to allocate any data structures to live across calls, you need to put them somewhere else. The memory context referenced by `multi_call_memory_ctx` is a suitable location for any data that needs to survive until the SRF is finished running. In most cases, this means that you should switch into `multi_call_memory_ctx` while doing the first-call setup. Use `funcctx->user_fctx` to hold a pointer to any such cross-call data structures. (Data you allocate in `multi_call_memory_ctx` will go away automatically when the query ends, so it is not necessary to free that data manually, either.)

Warning

While the actual arguments to the function remain unchanged between calls, if you detoast the argument values (which is normally done transparently by the `PG_GETARG_xxx` macro) in the transient context then the detoasted copies will be freed on each cycle. Accordingly, if you keep references to such values in your `user_fctx`, you must either copy them into the `multi_call_memory_ctx` after detoasting, or ensure that you detoast the values only in that context.

A complete pseudo-code example looks like the following:

```
Datum
my_set_returning_function(PG_FUNCTION_ARGS)
{
    FuncCallContext *funcctx;
    Datum          result;
    further declarations as needed

    if (SRF_IS_FIRSTCALL())
    {
        MemoryContext oldcontext;

        funcctx = SRF_FIRSTCALL_INIT();
        oldcontext = MemoryContextSwitchTo(funcctx->multi_call_memory_ctx);
        /* One-time setup code appears here: */
        user code
        if returning composite
            build TupleDesc, and perhaps AttInMetadata
        endif returning composite
        user code
        MemoryContextSwitchTo(oldcontext);
    }

    /* Each-time setup code appears here: */
    user code
    funcctx = SRF_PERCALL_SETUP();
    user code

    /* this is just one way we might test whether we are done: */
    if (funcctx->call_cntr < funcctx->max_calls)
    {
        /* Here we want to return another item: */
        user code
        obtain result Datum
        SRF_RETURN_NEXT(funcctx, result);
    }
    else
    {
        /* Here we are done returning items, so just report that fact. */
        /* (Resist the temptation to put cleanup code here.) */
    }
}
```

```

        SRF_RETURN_DONE(funcctx);
    }
}

```

A complete example of a simple SRF returning a composite type looks like:

```

PG_FUNCTION_INFO_V1(retcomposite);

Datum
retcomposite(PG_FUNCTION_ARGS)
{
    FuncCallContext    *funcctx;
    int                 call_cntr;
    int                 max_calls;
    TupleDesc           tupdesc;
    AttInMetadata       *attinmeta;

    /* stuff done only on the first call of the function */
    if (SRF_IS_FIRSTCALL())
    {
        MemoryContext    oldcontext;

        /* create a function context for cross-call persistence */
        funcctx = SRF_FIRSTCALL_INIT();

        /* switch to memory context appropriate for multiple function calls */
        oldcontext = MemoryContextSwitchTo(funcctx->multi_call_memory_ctx);

        /* total number of tuples to be returned */
        funcctx->max_calls = PG_GETARG_INT32(0);

        /* Build a tuple descriptor for our result type */
        if (get_call_result_type(fcinfo, NULL, &tupdesc) != TYPEFUNC_COMPOSITE)
            ereport(ERROR,
                    (errcode(ERRCODE_FEATURE_NOT_SUPPORTED),
                     errmsg("function returning record called in context "
                            "that cannot accept type record")));

        /*
         * generate attribute metadata needed later to produce tuples from raw
         * C strings
         */
        attinmeta = TupleDescGetAttInMetadata(tupdesc);
        funcctx->attinmeta = attinmeta;

        MemoryContextSwitchTo(oldcontext);
    }

    /* stuff done on every call of the function */
    funcctx = SRF_PERCALL_SETUP();

    call_cntr = funcctx->call_cntr;
    max_calls = funcctx->max_calls;
    attinmeta = funcctx->attinmeta;

    if (call_cntr < max_calls)    /* do when there is more left to send */
    {
        char            **values;
        HeapTuple        tuple;
    }
}

```

```

Datum            result;

/*
 * Prepare a values array for building the returned tuple.
 * This should be an array of C strings which will
 * be processed later by the type input functions.
 */
values = (char **) palloc(3 * sizeof(char *));
values[0] = (char *) palloc(16 * sizeof(char));
values[1] = (char *) palloc(16 * sizeof(char));
values[2] = (char *) palloc(16 * sizeof(char));

snprintf(values[0], 16, "%d", 1 * PG_GETARG_INT32(1));
snprintf(values[1], 16, "%d", 2 * PG_GETARG_INT32(1));
snprintf(values[2], 16, "%d", 3 * PG_GETARG_INT32(1));

/* build a tuple */
tuple = BuildTupleFromCStrings(atts, values);

/* make the tuple into a datum */
result = HeapTupleGetDatum(tuple);

/* clean up (this is not really necessary) */
pfree(values[0]);
pfree(values[1]);
pfree(values[2]);
pfree(values);

SRF_RETURN_NEXT(funcctx, result);
}
else /* do when there is no more left */
{
    SRF_RETURN_DONE(funcctx);
}
}

```

One way to declare this function in SQL is:

```

CREATE TYPE __retcomposite AS (f1 integer, f2 integer, f3 integer);

CREATE OR REPLACE FUNCTION retcomposite(integer, integer)
    RETURNS SETOF __retcomposite
    AS 'filename', 'retcomposite'
    LANGUAGE C IMMUTABLE STRICT;

```

A different way is to use OUT parameters:

```

CREATE OR REPLACE FUNCTION retcomposite(IN integer, IN integer,
    OUT f1 integer, OUT f2 integer, OUT f3 integer)
    RETURNS SETOF record
    AS 'filename', 'retcomposite'
    LANGUAGE C IMMUTABLE STRICT;

```

Notice that in this method the output type of the function is formally an anonymous `record` type.

39.10.9. Polymorphic Arguments and Return Types

C-language functions can be declared to accept and return the polymorphic types described in [Section 39.2.5](#). When a function's arguments or return types are defined as polymorphic types, the function author cannot know in advance what data type it will be called with, or need to return. There are two

routines provided in `fmgr.h` to allow a version-1 C function to discover the actual data types of its arguments and the type it is expected to return. The routines are called `get_fn_expr_rettype(FmgrInfo *flinfo)` and `get_fn_expr_argtype(FmgrInfo *flinfo, int argnum)`. They return the result or argument type OID, or `InvalidOid` if the information is not available. The structure `flinfo` is normally accessed as `fcinfo->flinfo`. The parameter `argnum` is zero based. `get_call_result_type` can also be used as an alternative to `get_fn_expr_rettype`. There is also `get_fn_expr_variadic`, which can be used to find out whether variadic arguments have been merged into an array. This is primarily useful for `VARIADIC "any"` functions, since such merging will always have occurred for variadic functions taking ordinary array types.

For example, suppose we want to write a function to accept a single element of any type, and return a one-dimensional array of that type:

```

PG_FUNCTION_INFO_V1(make_array);
Datum
make_array(PG_FUNCTION_ARGS)
{
    ArrayType *result;
    Oid        element_type = get_fn_expr_argtype(fcinfo->flinfo, 0);
    Datum      element;
    bool        isnull;
    int16       typelen;
    bool        typbyval;
    char        typalign;
    int         ndims;
    int         dims[MAXDIM];
    int         lbs[MAXDIM];

    if (!OidIsValid(element_type))
        elog(ERROR, "could not determine data type of input");

    /* get the provided element, being careful in case it's NULL */
    isnull = PG_ARGISNULL(0);
    if (isnull)
        element = (Datum) 0;
    else
        element = PG_GETARG_DATUM(0);

    /* we have one dimension */
    ndims = 1;
    /* and one element */
    dims[0] = 1;
    /* and lower bound is 1 */
    lbs[0] = 1;

    /* get required info about the element type */
    get_typlenbyvalalign(element_type, &typelen, &typbyval, &typalign);

    /* now build the array */
    result = construct_md_array(&element, &isnull, ndims, dims, lbs,
                               element_type, typelen, typbyval, typalign);

    PG_RETURN_ARRAYTYPE_P(result);
}

```

The following command declares the function `make_array` in SQL:

```

CREATE FUNCTION make_array(anyelement) RETURNS anyarray
AS 'DIRECTORY/funcs', 'make_array'

```

```
LANGUAGE C IMMUTABLE;
```

There is a variant of polymorphism that is only available to C-language functions: they can be declared to take parameters of type `"any"`. (Note that this type name must be double-quoted, since it's also an SQL reserved word.) This works like `anyelement` except that it does not constrain different `"any"` arguments to be the same type, nor do they help determine the function's result type. A C-language function can also declare its final parameter to be `VARIADIC "any"`. This will match one or more actual arguments of any type (not necessarily the same type). These arguments will *not* be gathered into an array as happens with normal variadic functions; they will just be passed to the function separately. The `PG_NARGS()` macro and the methods described above must be used to determine the number of actual arguments and their types when using this feature. Also, users of such a function might wish to use the `VARIADIC` keyword in their function call, with the expectation that the function would treat the array elements as separate arguments. The function itself must implement that behavior if wanted, after using `get_fn_expr_variadic` to detect that the actual argument was marked with `VARIADIC`.

39.10.10. Shared Memory

39.10.10.1. Requesting Shared Memory at Startup

Add-ins can reserve shared memory on server startup. To do so, the add-in's shared library must be preloaded by specifying it in [shared_preload_libraries](#). The shared library should also register a `shmem_request_hook` in its `_PG_init` function. This `shmem_request_hook` can reserve shared memory by calling:

```
void RequestAddinShmemSpace(Size size)
```

Each backend should obtain a pointer to the reserved shared memory by calling:

```
void *ShmemInitStruct(const char *name, Size size, bool *foundPtr)
```

If this function sets `foundPtr` to `false`, the caller should proceed to initialize the contents of the reserved shared memory. If `foundPtr` is set to `true`, the shared memory was already initialized by another backend, and the caller need not initialize further.

To avoid race conditions, each backend should use the `LWLock AddinShmemInitLock` when initializing its allocation of shared memory, as shown here:

```
static mystruct *ptr = NULL;
bool          found;

LWLockAcquire(AddinShmemInitLock, LW_EXCLUSIVE);
ptr = ShmemInitStruct("my struct name", size, &found);
if (!found)
{
    ... initialize contents of shared memory ...
    ptr->locks = GetNamedLWLockTranche("my tranche name");
}
LWLockRelease(AddinShmemInitLock);
```

`shmem_startup_hook` provides a convenient place for the initialization code, but it is not strictly required that all such code be placed in this hook. Each backend will execute the registered `shmem_startup_hook` shortly after it attaches to shared memory. Note that add-ins should still acquire `AddinShmemInitLock` within this hook, as shown in the example above.

39.10.10.2. Requesting Shared Memory After Startup

There is another, more flexible method of reserving shared memory that can be done after server startup and outside a `shmem_request_hook`. To do so, each backend that will use the shared memory should obtain a pointer to it by calling:

```
void *GetNamedDSMSegment(const char *name, size_t size,
                        void (*init_callback) (void *ptr),
```

```
bool *found)
```

If a dynamic shared memory segment with the given name does not yet exist, this function will allocate it and initialize it with the provided `init_callback` callback function. If the segment has already been allocated and initialized by another backend, this function simply attaches the existing dynamic shared memory segment to the current backend.

Unlike shared memory reserved at server startup, there is no need to acquire `AddinShmemInitLock` or otherwise take action to avoid race conditions when reserving shared memory with `GetNamedDSMSegment`. This function ensures that only one backend allocates and initializes the segment and that all other backends receive a pointer to the fully allocated and initialized segment.

39.10.11. LWLocks

39.10.11.1. Requesting LWLocks at Startup

Add-ins can reserve LWLocks on server startup. As with shared memory reserved at server startup, the add-in's shared library must be preloaded by specifying it in [shared preload libraries](#), and the shared library should register a `shmem_request_hook` in its `_PG_init` function. This `shmem_request_hook` can reserve LWLocks by calling:

```
void RequestNamedLWLockTranche(const char *tranche_name, int num_lwlocks)
```

This ensures that an array of `num_lwlocks` LWLocks is available under the name `tranche_name`. A pointer to this array can be obtained by calling:

```
LWLockPadded *GetNamedLWLockTranche(const char *tranche_name)
```

39.10.11.2. Requesting LWLocks After Startup

There is another, more flexible method of obtaining LWLocks that can be done after server startup and outside a `shmem_request_hook`. To do so, first allocate a `tranche_id` by calling:

```
int LWLockNewTrancheId(void)
```

Next, initialize each LWLock, passing the new `tranche_id` as an argument:

```
void LWLockInitialize(LWLock *lock, int tranche_id)
```

Similar to shared memory, each backend should ensure that only one process allocates a new `tranche_id` and initializes each new LWLock. One way to do this is to only call these functions in your shared memory initialization code with the `AddinShmemInitLock` held exclusively. If using `GetNamedDSMSegment`, calling these functions in the `init_callback` callback function is sufficient to avoid race conditions.

Finally, each backend using the `tranche_id` should associate it with a `tranche_name` by calling:

```
void LWLockRegisterTranche(int tranche_id, const char *tranche_name)
```

39.10.12. Custom Wait Events

Add-ins can define custom wait events under the wait event type `Extension` by calling:

```
uint32 WaitEventExtensionNew(const char *wait_event_name)
```

The wait event is associated to a user-facing custom string.

Custom wait events can be viewed in [pg_stat_activity](#):

```
=# SELECT wait_event_type, wait_event FROM pg_stat_activity
   WHERE backend_type ~ 'worker_spi';
 wait_event_type | wait_event
-----+-----
 Extension      | WorkerSpiMain
```


(1 row)

39.10.13. Injection Points

An injection point with a given `name` is declared using macro:

```
INJECTION_POINT(name);
```

There are a few injection points already declared at strategic points within the server code. After adding a new injection point the code needs to be compiled in order for that injection point to be available in the binary. Add-ins written in C-language can declare injection points in their own code using the same macro. The injection point names should use lower-case characters, with terms separated by dashes.

Add-ins can attach callbacks to an already-declared injection point by calling:

```
extern void InjectionPointAttach(const char *name,
                                const char *library,
                                const char *function,
                                const void *private_data,
                                int private_data_size);
```

`name` is the name of the injection point, which when reached during execution will execute the function loaded from `library`. `private_data` is a private area of data of size `private_data_size` given as argument to the callback when executed.

Here is an example of callback for `InjectionPointCallback`:

```
static void
custom_injection_callback(const char *name, const void *private_data)
{
    uint32 wait_event_info = WaitEventInjectionPointNew(name);

    pgstat_report_wait_start(wait_event_info);
    elog(NOTICE, "%s: executed custom callback", name);
    pgstat_report_wait_end();
}
```

This callback prints a message to server error log with severity `NOTICE`, but callbacks may implement more complex logic.

Optionally, it is possible to detach an injection point by calling:

```
extern bool InjectionPointDetach(const char *name);
```

On success, `true` is returned, `false` otherwise.

A callback attached to an injection point is available across all the backends including the backends started after `InjectionPointAttach` is called. It remains attached while the server is running or until the injection point is detached using `InjectionPointDetach`.

Injection points are only available if Postgres Pro is compiled with support for it.

39.10.14. Using C++ for Extensibility

Although the Postgres Pro backend is written in C, it is possible to write extensions in C++ if these guidelines are followed:

- All functions accessed by the backend must present a C interface to the backend; these C functions can then call C++ functions. For example, `extern C` linkage is required for backend-accessed functions. This is also necessary for any functions that are passed as pointers between the backend and C++ code.
- Free memory using the appropriate deallocation method. For example, most backend memory is allocated using `malloc()`, so use `free()` to free it. Using C++ `delete` in such cases will fail.

- Prevent exceptions from propagating into the C code (use a catch-all block at the top level of all `extern C` functions). This is necessary even if the C++ code does not explicitly throw any exceptions, because events like out-of-memory can still throw exceptions. Any exceptions must be caught and appropriate errors passed back to the C interface. If possible, compile C++ with `-fno-exceptions` to eliminate exceptions entirely; in such cases, you must check for failures in your C++ code, e.g., check for NULL returned by `new()`.
- If calling backend functions from C++ code, be sure that the C++ call stack contains only plain old data structures (POD). This is necessary because backend errors generate a distant `longjmp()` that does not properly unroll a C++ call stack with non-POD objects.

In summary, it is best to place C++ code behind a wall of `extern C` functions that interface to the backend, and avoid exception, memory, and call stack leakage.

39.11. Function Optimization Information

By default, a function is just a “black box” that the database system knows very little about the behavior of. However, that means that queries using the function may be executed much less efficiently than they could be. It is possible to supply additional knowledge that helps the planner optimize function calls.

Some basic facts can be supplied by declarative annotations provided in the `CREATE FUNCTION` command. Most important of these is the function's **volatility category** (`IMMUTABLE`, `STABLE`, or `VOLATILE`); one should always be careful to specify this correctly when defining a function. The parallel safety property (`PARALLEL UNSAFE`, `PARALLEL RESTRICTED`, or `PARALLEL SAFE`) must also be specified if you hope to use the function in parallelized queries. It can also be useful to specify the function's estimated execution cost, and/or the number of rows a set-returning function is estimated to return. However, the declarative way of specifying those two facts only allows specifying a constant value, which is often inadequate.

It is also possible to attach a *planner support function* to an SQL-callable function (called its *target function*), and thereby provide knowledge about the target function that is too complex to be represented declaratively. Planner support functions have to be written in C (although their target functions might not be), so this is an advanced feature that relatively few people will use.

A planner support function must have the SQL signature

```
supportfn(internal) returns internal
```

It is attached to its target function by specifying the `SUPPORT` clause when creating the target function.

Here we provide an overview of what planner support functions can do. The set of possible requests to a support function is extensible, so more things might be possible in future versions.

Some function calls can be simplified during planning based on properties specific to the function. For example, `int4mul(n, 1)` could be simplified to just `n`. This type of transformation can be performed by a planner support function, by having it implement the `SupportRequestSimplify` request type. The support function will be called for each instance of its target function found in a query parse tree. If it finds that the particular call can be simplified into some other form, it can build and return a parse tree representing that expression. This will automatically work for operators based on the function, too — in the example just given, `n * 1` would also be simplified to `n`. (But note that this is just an example; this particular optimization is not actually performed by standard Postgres Pro.) We make no guarantee that Postgres Pro will never call the target function in cases that the support function could simplify. Ensure rigorous equivalence between the simplified expression and an actual execution of the target function.

For target functions that return `boolean`, it is often useful to estimate the fraction of rows that will be selected by a `WHERE` clause using that function. This can be done by a support function that implements the `SupportRequestSelectivity` request type.

If the target function's run time is highly dependent on its inputs, it may be useful to provide a non-constant cost estimate for it. This can be done by a support function that implements the `SupportRequestCost` request type.

For target functions that return sets, it is often useful to provide a non-constant estimate for the number of rows that will be returned. This can be done by a support function that implements the `SupportRequestRows` request type.

For target functions that return `boolean`, it may be possible to convert a function call appearing in `WHERE` into an indexable operator clause or clauses. The converted clauses might be exactly equivalent to the function's condition, or they could be somewhat weaker (that is, they might accept some values that the function condition does not). In the latter case the index condition is said to be *lossy*; it can still be used to scan an index, but the function call will have to be executed for each row returned by the index to see if it really passes the `WHERE` condition or not. To create such conditions, the support function must implement the `SupportRequestIndexCondition` request type.

39.12. User-Defined Aggregates

Aggregate functions in Postgres Pro are defined in terms of *state values* and *state transition functions*. That is, an aggregate operates using a state value that is updated as each successive input row is processed. To define a new aggregate function, one selects a data type for the state value, an initial value for the state, and a state transition function. The state transition function takes the previous state value and the aggregate's input value(s) for the current row, and returns a new state value. A *final function* can also be specified, in case the desired result of the aggregate is different from the data that needs to be kept in the running state value. The final function takes the ending state value and returns whatever is wanted as the aggregate result. In principle, the transition and final functions are just ordinary functions that could also be used outside the context of the aggregate. (In practice, it's often helpful for performance reasons to create specialized transition functions that can only work when called as part of an aggregate.)

Thus, in addition to the argument and result data types seen by a user of the aggregate, there is an internal state-value data type that might be different from both the argument and result types.

If we define an aggregate that does not use a final function, we have an aggregate that computes a running function of the column values from each row. `sum` is an example of this kind of aggregate. `sum` starts at zero and always adds the current row's value to its running total. For example, if we want to make a `sum` aggregate to work on a data type for complex numbers, we only need the addition function for that data type. The aggregate definition would be:

```
CREATE AGGREGATE sum (complex)
(
    sfunc = complex_add,
    stype = complex,
    initcond = '(0,0)'
);
```

which we might use like this:

```
SELECT sum(a) FROM test_complex;
```

```
sum
-----
(34,53.9)
```

(Notice that we are relying on function overloading: there is more than one aggregate named `sum`, but Postgres Pro can figure out which kind of `sum` applies to a column of type `complex`.)

The above definition of `sum` will return zero (the initial state value) if there are no nonnull input values. Perhaps we want to return null in that case instead — the SQL standard expects `sum` to behave that way. We can do this simply by omitting the `initcond` phrase, so that the initial state value is null. Ordinarily this would mean that the `sfunc` would need to check for a null state-value input. But for `sum` and some other simple aggregates like `max` and `min`, it is sufficient to insert the first nonnull input value into the state variable and then start applying the transition function at the second nonnull input value. Postgres Pro will do that automatically if the initial state value is null and the transition function is marked “strict” (i.e., not to be called for null inputs).

Another bit of default behavior for a “strict” transition function is that the previous state value is retained unchanged whenever a null input value is encountered. Thus, null values are ignored. If you need some other behavior for null inputs, do not declare your transition function as strict; instead code it to test for null inputs and do whatever is needed.

`avg` (average) is a more complex example of an aggregate. It requires two pieces of running state: the sum of the inputs and the count of the number of inputs. The final result is obtained by dividing these quantities. Average is typically implemented by using an array as the state value. For example, the built-in implementation of `avg(float8)` looks like:

```
CREATE AGGREGATE avg (float8)
(
    sfunc = float8_accum,
    stype = float8[],
    finalfunc = float8_avg,
    initcond = '{0,0,0}'
);
```

Note

`float8_accum` requires a three-element array, not just two elements, because it accumulates the sum of squares as well as the sum and count of the inputs. This is so that it can be used for some other aggregates as well as `avg`.

Aggregate function calls in SQL allow `DISTINCT` and `ORDER BY` options that control which rows are fed to the aggregate's transition function and in what order. These options are implemented behind the scenes and are not the concern of the aggregate's support functions.

For further details see the [CREATE AGGREGATE](#) command.

39.12.1. Moving-Aggregate Mode

Aggregate functions can optionally support *moving-aggregate mode*, which allows substantially faster execution of aggregate functions within windows with moving frame starting points. (See [Section 3.5](#) and [Section 4.2.8](#) for information about use of aggregate functions as window functions.) The basic idea is that in addition to a normal “forward” transition function, the aggregate provides an *inverse transition function*, which allows rows to be removed from the aggregate's running state value when they exit the window frame. For example a `sum` aggregate, which uses addition as the forward transition function, would use subtraction as the inverse transition function. Without an inverse transition function, the window function mechanism must recalculate the aggregate from scratch each time the frame starting point moves, resulting in run time proportional to the number of input rows times the average frame length. With an inverse transition function, the run time is only proportional to the number of input rows.

The inverse transition function is passed the current state value and the aggregate input value(s) for the earliest row included in the current state. It must reconstruct what the state value would have been if the given input row had never been aggregated, but only the rows following it. This sometimes requires that the forward transition function keep more state than is needed for plain aggregation mode. Therefore, the moving-aggregate mode uses a completely separate implementation from the plain mode: it has its own state data type, its own forward transition function, and its own final function if needed. These can be the same as the plain mode's data type and functions, if there is no need for extra state.

As an example, we could extend the `sum` aggregate given above to support moving-aggregate mode like this:

```
CREATE AGGREGATE sum (complex)
(
    sfunc = complex_add,
    stype = complex,
```

```
initcond = '(0,0)',
msfunc = complex_add,
minvfunc = complex_sub,
mstype = complex,
minitcond = '(0,0)'
);
```

The parameters whose names begin with `m` define the moving-aggregate implementation. Except for the inverse transition function `minvfunc`, they correspond to the plain-aggregate parameters without `m`.

The forward transition function for moving-aggregate mode is not allowed to return null as the new state value. If the inverse transition function returns null, this is taken as an indication that the inverse function cannot reverse the state calculation for this particular input, and so the aggregate calculation will be redone from scratch for the current frame starting position. This convention allows moving-aggregate mode to be used in situations where there are some infrequent cases that are impractical to reverse out of the running state value. The inverse transition function can “punt” on these cases, and yet still come out ahead so long as it can work for most cases. As an example, an aggregate working with floating-point numbers might choose to punt when a NaN (not a number) input has to be removed from the running state value.

When writing moving-aggregate support functions, it is important to be sure that the inverse transition function can reconstruct the correct state value exactly. Otherwise there might be user-visible differences in results depending on whether the moving-aggregate mode is used. An example of an aggregate for which adding an inverse transition function seems easy at first, yet where this requirement cannot be met is `sum` over `float4` or `float8` inputs. A naive declaration of `sum(float8)` could be

```
CREATE AGGREGATE unsafe_sum (float8)
(
    stype = float8,
    sfunc = float8pl,
    mstype = float8,
    msfunc = float8pl,
    minvfunc = float8mi
);
```

This aggregate, however, can give wildly different results than it would have without the inverse transition function. For example, consider

```
SELECT
    unsafe_sum(x) OVER (ORDER BY n ROWS BETWEEN CURRENT ROW AND 1 FOLLOWING)
FROM (VALUES (1, 1.0e20::float8),
            (2, 1.0::float8)) AS v (n,x);
```

This query returns 0 as its second result, rather than the expected answer of 1. The cause is the limited precision of floating-point values: adding 1 to `1e20` results in `1e20` again, and so subtracting `1e20` from that yields 0, not 1. Note that this is a limitation of floating-point arithmetic in general, not a limitation of Postgres Pro.

39.12.2. Polymorphic and Variadic Aggregates

Aggregate functions can use polymorphic state transition functions or final functions, so that the same functions can be used to implement multiple aggregates. See [Section 39.2.5](#) for an explanation of polymorphic functions. Going a step further, the aggregate function itself can be specified with polymorphic input type(s) and state type, allowing a single aggregate definition to serve for multiple input data types. Here is an example of a polymorphic aggregate:

```
CREATE AGGREGATE array_accum (anycompatible)
(
    sfunc = array_append,
    stype = anycompatiblearray,
    initcond = '{}'
```

```
);
```

Here, the actual state type for any given aggregate call is the array type having the actual input type as elements. The behavior of the aggregate is to concatenate all the inputs into an array of that type. (Note: the built-in aggregate `array_agg` provides similar functionality, with better performance than this definition would have.)

Here's the output using two different actual data types as arguments:

```
SELECT attrelid::regclass, array_accum(attname)
FROM pg_attribute
WHERE attnum > 0 AND attrelid = 'pg_tablespace'::regclass
GROUP BY attrelid;
```

attrelid	array_accum
pg_tablespace	{spcname, spcowner, spcacl, spcoptions}

(1 row)

```
SELECT attrelid::regclass, array_accum(atttypid::regtype)
FROM pg_attribute
WHERE attnum > 0 AND attrelid = 'pg_tablespace'::regclass
GROUP BY attrelid;
```

attrelid	array_accum
pg_tablespace	{name, oid, aclitem[], text[]}

(1 row)

Ordinarily, an aggregate function with a polymorphic result type has a polymorphic state type, as in the above example. This is necessary because otherwise the final function cannot be declared sensibly: it would need to have a polymorphic result type but no polymorphic argument type, which `CREATE FUNCTION` will reject on the grounds that the result type cannot be deduced from a call. But sometimes it is inconvenient to use a polymorphic state type. The most common case is where the aggregate support functions are to be written in C and the state type should be declared as `internal` because there is no SQL-level equivalent for it. To address this case, it is possible to declare the final function as taking extra “dummy” arguments that match the input arguments of the aggregate. Such dummy arguments are always passed as null values since no specific value is available when the final function is called. Their only use is to allow a polymorphic final function's result type to be connected to the aggregate's input type(s). For example, the definition of the built-in aggregate `array_agg` is equivalent to

```
CREATE FUNCTION array_agg_transfn(internal, anynonarray)
RETURNS internal ...;
CREATE FUNCTION array_agg_finalfn(internal, anynonarray)
RETURNS anyarray ...;

CREATE AGGREGATE array_agg (anynonarray)
(
    sfunc = array_agg_transfn,
    stype = internal,
    finalfunc = array_agg_finalfn,
    finalfunc_extra
);
```

Here, the `finalfunc_extra` option specifies that the final function receives, in addition to the state value, extra dummy argument(s) corresponding to the aggregate's input argument(s). The extra `anynonarray` argument allows the declaration of `array_agg_finalfn` to be valid.

An aggregate function can be made to accept a varying number of arguments by declaring its last argument as a `VARIADIC` array, in much the same fashion as for regular functions; see [Section 39.5.6](#). The

aggregate's transition function(s) must have the same array type as their last argument. The transition function(s) typically would also be marked `VARIADIC`, but this is not strictly required.

Note

Variadic aggregates are easily misused in connection with the `ORDER BY` option (see [Section 4.2.7](#)), since the parser cannot tell whether the wrong number of actual arguments have been given in such a combination. Keep in mind that everything to the right of `ORDER BY` is a sort key, not an argument to the aggregate. For example, in

```
SELECT myaggregate(a ORDER BY a, b, c) FROM ...
```

the parser will see this as a single aggregate function argument and three sort keys. However, the user might have intended

```
SELECT myaggregate(a, b, c ORDER BY a) FROM ...
```

If `myaggregate` is variadic, both these calls could be perfectly valid.

For the same reason, it's wise to think twice before creating aggregate functions with the same names and different numbers of regular arguments.

39.12.3. Ordered-Set Aggregates

The aggregates we have been describing so far are “normal” aggregates. Postgres Pro also supports *ordered-set aggregates*, which differ from normal aggregates in two key ways. First, in addition to ordinary aggregated arguments that are evaluated once per input row, an ordered-set aggregate can have “direct” arguments that are evaluated only once per aggregation operation. Second, the syntax for the ordinary aggregated arguments specifies a sort ordering for them explicitly. An ordered-set aggregate is usually used to implement a computation that depends on a specific row ordering, for instance rank or percentile, so that the sort ordering is a required aspect of any call. For example, the built-in definition of `percentile_disc` is equivalent to:

```
CREATE FUNCTION ordered_set_transition(internal, anyelement)
  RETURNS internal ...;
CREATE FUNCTION percentile_disc_final(internal, float8, anyelement)
  RETURNS anyelement ...;

CREATE AGGREGATE percentile_disc (float8 ORDER BY anyelement)
(
  sfunc = ordered_set_transition,
  stype = internal,
  finalfunc = percentile_disc_final,
  finalfunc_extra
);
```

This aggregate takes a `float8` direct argument (the percentile fraction) and an aggregated input that can be of any sortable data type. It could be used to obtain a median household income like this:

```
SELECT percentile_disc(0.5) WITHIN GROUP (ORDER BY income) FROM households;
percentile_disc
-----
50489
```

Here, `0.5` is a direct argument; it would make no sense for the percentile fraction to be a value varying across rows.

Unlike the case for normal aggregates, the sorting of input rows for an ordered-set aggregate is *not* done behind the scenes, but is the responsibility of the aggregate's support functions. The typical implementation approach is to keep a reference to a “tuplesort” object in the aggregate's state value, feed

the incoming rows into that object, and then complete the sorting and read out the data in the final function. This design allows the final function to perform special operations such as injecting additional “hypothetical” rows into the data to be sorted. While normal aggregates can often be implemented with support functions written in PL/pgSQL or another PL language, ordered-set aggregates generally have to be written in C, since their state values aren't definable as any SQL data type. (In the above example, notice that the state value is declared as type `internal` — this is typical.) Also, because the final function performs the sort, it is not possible to continue adding input rows by executing the transition function again later. This means the final function is not `READ_ONLY`; it must be declared in `CREATE AGGREGATE` as `READ_WRITE`, or as `SHAREABLE` if it's possible for additional final-function calls to make use of the already-sorted state.

The state transition function for an ordered-set aggregate receives the current state value plus the aggregated input values for each row, and returns the updated state value. This is the same definition as for normal aggregates, but note that the direct arguments (if any) are not provided. The final function receives the last state value, the values of the direct arguments if any, and (if `finalfunc_extra` is specified) null values corresponding to the aggregated input(s). As with normal aggregates, `finalfunc_extra` is only really useful if the aggregate is polymorphic; then the extra dummy argument(s) are needed to connect the final function's result type to the aggregate's input type(s).

Currently, ordered-set aggregates cannot be used as window functions, and therefore there is no need for them to support moving-aggregate mode.

39.12.4. Partial Aggregation

Optionally, an aggregate function can support *partial aggregation*. The idea of partial aggregation is to run the aggregate's state transition function over different subsets of the input data independently, and then to combine the state values resulting from those subsets to produce the same state value that would have resulted from scanning all the input in a single operation. This mode can be used for parallel aggregation by having different worker processes scan different portions of a table. Each worker produces a partial state value, and at the end those state values are combined to produce a final state value. (In the future this mode might also be used for purposes such as combining aggregations over local and remote tables; but that is not implemented yet.)

To support partial aggregation, the aggregate definition must provide a *combine function*, which takes two values of the aggregate's state type (representing the results of aggregating over two subsets of the input rows) and produces a new value of the state type, representing what the state would have been after aggregating over the combination of those sets of rows. It is unspecified what the relative order of the input rows from the two sets would have been. This means that it's usually impossible to define a useful combine function for aggregates that are sensitive to input row order.

As simple examples, `MAX` and `MIN` aggregates can be made to support partial aggregation by specifying the combine function as the same greater-of-two or lesser-of-two comparison function that is used as their transition function. `SUM` aggregates just need an addition function as combine function. (Again, this is the same as their transition function, unless the state value is wider than the input data type.)

The combine function is treated much like a transition function that happens to take a value of the state type, not of the underlying input type, as its second argument. In particular, the rules for dealing with null values and strict functions are similar. Also, if the aggregate definition specifies a non-null `initcond`, keep in mind that that will be used not only as the initial state for each partial aggregation run, but also as the initial state for the combine function, which will be called to combine each partial result into that state.

If the aggregate's state type is declared as `internal`, it is the combine function's responsibility that its result is allocated in the correct memory context for aggregate state values. This means in particular that when the first input is `NULL` it's invalid to simply return the second input, as that value will be in the wrong context and will not have sufficient lifespan.

When the aggregate's state type is declared as `internal`, it is usually also appropriate for the aggregate definition to provide a *serialization function* and a *deserialization function*, which allow such a state

value to be copied from one process to another. Without these functions, parallel aggregation cannot be performed, and future applications such as local/remote aggregation will probably not work either.

A serialization function must take a single argument of type `internal` and return a result of type `bytea`, which represents the state value packaged up into a flat blob of bytes. Conversely, a deserialization function reverses that conversion. It must take two arguments of types `bytea` and `internal`, and return a result of type `internal`. (The second argument is unused and is always zero, but it is required for type-safety reasons.) The result of the deserialization function should simply be allocated in the current memory context, as unlike the combine function's result, it is not long-lived.

Worth noting also is that for an aggregate to be executed in parallel, the aggregate itself must be marked `PARALLEL SAFE`. The parallel-safety markings on its support functions are not consulted.

39.12.5. Support Functions for Aggregates

A function written in C can detect that it is being called as an aggregate support function by calling `AggCheckCallContext`, for example:

```
if (AggCheckCallContext(fcinfo, NULL))
```

One reason for checking this is that when it is true, the first input must be a temporary state value and can therefore safely be modified in-place rather than allocating a new copy. See `int8inc()` for an example. (While aggregate transition functions are always allowed to modify the transition value in-place, aggregate final functions are generally discouraged from doing so; if they do so, the behavior must be declared when creating the aggregate. See [CREATE AGGREGATE](#) for more detail.)

The second argument of `AggCheckCallContext` can be used to retrieve the memory context in which aggregate state values are being kept. This is useful for transition functions that wish to use “expanded” objects (see [Section 39.13.1](#)) as their state values. On first call, the transition function should return an expanded object whose memory context is a child of the aggregate state context, and then keep returning the same expanded object on subsequent calls. See `array_append()` for an example. (`array_append()` is not the transition function of any built-in aggregate, but it is written to behave efficiently when used as transition function of a custom aggregate.)

Another support routine available to aggregate functions written in C is `AggGetAggref`, which returns the `Aggref` parse node that defines the aggregate call. This is mainly useful for ordered-set aggregates, which can inspect the substructure of the `Aggref` node to find out what sort ordering they are supposed to implement.

39.13. User-Defined Types

As described in [Section 39.2](#), Postgres Pro can be extended to support new data types. This section describes how to define new base types, which are data types defined below the level of the SQL language. Creating a new base type requires implementing functions to operate on the type in a low-level language, usually C.

A user-defined type must always have input and output functions. These functions determine how the type appears in strings (for input by the user and output to the user) and how the type is organized in memory. The input function takes a null-terminated character string as its argument and returns the internal (in memory) representation of the type. The output function takes the internal representation of the type as argument and returns a null-terminated character string. If we want to do anything more with the type than merely store it, we must provide additional functions to implement whatever operations we'd like to have for the type.

Suppose we want to define a type `complex` that represents complex numbers. A natural way to represent a complex number in memory would be the following C structure:

```
typedef struct Complex {
    double    x;
    double    y;
```

```
} Complex;
```

We will need to make this a pass-by-reference type, since it's too large to fit into a single `Datum` value.

As the external string representation of the type, we choose a string of the form (x,y) .

The input and output functions are usually not hard to write, especially the output function. But when defining the external string representation of the type, remember that you must eventually write a complete and robust parser for that representation as your input function. For instance:

```
PG_FUNCTION_INFO_V1(complex_in);

Datum
complex_in(PG_FUNCTION_ARGS)
{
    char          *str = PG_GETARG_CSTRING(0);
    double        x,
                  y;
    Complex       *result;

    if (sscanf(str, " ( %lf , %lf )", &x, &y) != 2)
        ereport(ERROR,
                (errcode(ERRCODE_INVALID_TEXT_REPRESENTATION),
                 errmsg("invalid input syntax for type %s: \"%s\"",
                        "complex", str)));

    result = (Complex *) palloc(sizeof(Complex));
    result->x = x;
    result->y = y;
    PG_RETURN_POINTER(result);
}
```

The output function can simply be:

```
PG_FUNCTION_INFO_V1(complex_out);

Datum
complex_out(PG_FUNCTION_ARGS)
{
    Complex       *complex = (Complex *) PG_GETARG_POINTER(0);
    char          *result;

    result = psprintf("(%g,%g)", complex->x, complex->y);
    PG_RETURN_CSTRING(result);
}
```

You should be careful to make the input and output functions inverses of each other. If you do not, you will have severe problems when you need to dump your data into a file and then read it back in. This is a particularly common problem when floating-point numbers are involved.

Optionally, a user-defined type can provide binary input and output routines. Binary I/O is normally faster but less portable than textual I/O. As with textual I/O, it is up to you to define exactly what the external binary representation is. Most of the built-in data types try to provide a machine-independent binary representation. For `complex`, we will piggy-back on the binary I/O converters for type `float8`:

```
PG_FUNCTION_INFO_V1(complex_recv);

Datum
complex_recv(PG_FUNCTION_ARGS)
```

```

{
    StringInfo  buf = (StringInfo) PG_GETARG_POINTER(0);
    Complex    *result;

    result = (Complex *) palloc(sizeof(Complex));
    result->x = pq_getmsgfloat8(buf);
    result->y = pq_getmsgfloat8(buf);
    PG_RETURN_POINTER(result);
}

PG_FUNCTION_INFO_V1(complex_send);

Datum
complex_send(PG_FUNCTION_ARGS)
{
    Complex    *complex = (Complex *) PG_GETARG_POINTER(0);
    StringInfoData buf;

    pq_begintypsend(&buf);
    pq_sendfloat8(&buf, complex->x);
    pq_sendfloat8(&buf, complex->y);
    PG_RETURN_BYTEA_P(pq_endtypsend(&buf));
}

```

Once we have written the I/O functions and compiled them into a shared library, we can define the `complex` type in SQL. First we declare it as a shell type:

```
CREATE TYPE complex;
```

This serves as a placeholder that allows us to reference the type while defining its I/O functions. Now we can define the I/O functions:

```

CREATE FUNCTION complex_in(cstring)
    RETURNS complex
    AS 'filename'
    LANGUAGE C IMMUTABLE STRICT;

CREATE FUNCTION complex_out(complex)
    RETURNS cstring
    AS 'filename'
    LANGUAGE C IMMUTABLE STRICT;

CREATE FUNCTION complex_recv(internal)
    RETURNS complex
    AS 'filename'
    LANGUAGE C IMMUTABLE STRICT;

CREATE FUNCTION complex_send(complex)
    RETURNS bytea
    AS 'filename'
    LANGUAGE C IMMUTABLE STRICT;

```

Finally, we can provide the full definition of the data type:

```

CREATE TYPE complex (
    internallength = 16,
    input = complex_in,
    output = complex_out,
    receive = complex_recv,

```

```

    send = complex_send,
    alignment = double
);

```

When you define a new base type, Postgres Pro automatically provides support for arrays of that type. The array type typically has the same name as the base type with the underscore character (`_`) prepended.

Once the data type exists, we can declare additional functions to provide useful operations on the data type. Operators can then be defined atop the functions, and if needed, operator classes can be created to support indexing of the data type. These additional layers are discussed in following sections.

If the internal representation of the data type is variable-length, the internal representation must follow the standard layout for variable-length data: the first four bytes must be a `char[4]` field which is never accessed directly (customarily named `vl_len_`). You must use the `SET_VARSIZE()` macro to store the total size of the datum (including the length field itself) in this field and `VARSIZE()` to retrieve it. (These macros exist because the length field may be encoded depending on platform.)

For further details see the description of the [CREATE TYPE](#) command.

39.13.1. TOAST Considerations

If the values of your data type vary in size (in internal form), it's usually desirable to make the data type TOAST-able (see [Section 66.2](#)). You should do this even if the values are always too small to be compressed or stored externally, because TOAST can save space on small data too, by reducing header overhead.

To support TOAST storage, the C functions operating on the data type must always be careful to unpack any toasted values they are handed by using `PG_DETOAST_DATUM`. (This detail is customarily hidden by defining type-specific `GETARG_DATATYPE_P` macros.) Then, when running the `CREATE TYPE` command, specify the internal length as `variable` and select some appropriate storage option other than `plain`.

If data alignment is unimportant (either just for a specific function or because the data type specifies byte alignment anyway) then it's possible to avoid some of the overhead of `PG_DETOAST_DATUM`. You can use `PG_DETOAST_DATUM_PACKED` instead (customarily hidden by defining a `GETARG_DATATYPE_PP` macro) and using the macros `VARSIZE_ANY_EXHDR` and `VARDATA_ANY` to access a potentially-packed datum. Again, the data returned by these macros is not aligned even if the data type definition specifies an alignment. If the alignment is important you must go through the regular `PG_DETOAST_DATUM` interface.

Note

Older code frequently declares `vl_len_` as an `int32` field instead of `char[4]`. This is OK as long as the struct definition has other fields that have at least `int32` alignment. But it is dangerous to use such a struct definition when working with a potentially unaligned datum; the compiler may take it as license to assume the datum actually is aligned, leading to core dumps on architectures that are strict about alignment.

Another feature that's enabled by TOAST support is the possibility of having an *expanded* in-memory data representation that is more convenient to work with than the format that is stored on disk. The regular or “flat” varlena storage format is ultimately just a blob of bytes; it cannot for example contain pointers, since it may get copied to other locations in memory. For complex data types, the flat format may be quite expensive to work with, so Postgres Pro provides a way to “expand” the flat format into a representation that is more suited to computation, and then pass that format in-memory between functions of the data type.

To use expanded storage, a data type must define an expanded format that follows the rules given in [src/include/utils/expandeddatum.h](#), and provide functions to “expand” a flat varlena value into expanded format and “flatten” the expanded format back to the regular varlena representation. Then ensure that

all C functions for the data type can accept either representation, possibly by converting one into the other immediately upon receipt. This does not require fixing all existing functions for the data type at once, because the standard `PG_DETOAST_DATUM` macro is defined to convert expanded inputs into regular flat format. Therefore, existing functions that work with the flat varlena format will continue to work, though slightly inefficiently, with expanded inputs; they need not be converted until and unless better performance is important.

C functions that know how to work with an expanded representation typically fall into two categories: those that can only handle expanded format, and those that can handle either expanded or flat varlena inputs. The former are easier to write but may be less efficient overall, because converting a flat input to expanded form for use by a single function may cost more than is saved by operating on the expanded format. When only expanded format need be handled, conversion of flat inputs to expanded form can be hidden inside an argument-fetching macro, so that the function appears no more complex than one working with traditional varlena input. To handle both types of input, write an argument-fetching function that will detoast external, short-header, and compressed varlena inputs, but not expanded inputs. Such a function can be defined as returning a pointer to a union of the flat varlena format and the expanded format. Callers can use the `VARATT_IS_EXPANDED_HEADER()` macro to determine which format they received.

The TOAST infrastructure not only allows regular varlena values to be distinguished from expanded values, but also distinguishes “read-write” and “read-only” pointers to expanded values. C functions that only need to examine an expanded value, or will only change it in safe and non-semantically-visible ways, need not care which type of pointer they receive. C functions that produce a modified version of an input value are allowed to modify an expanded input value in-place if they receive a read-write pointer, but must not modify the input if they receive a read-only pointer; in that case they have to copy the value first, producing a new value to modify. A C function that has constructed a new expanded value should always return a read-write pointer to it. Also, a C function that is modifying a read-write expanded value in-place should take care to leave the value in a sane state if it fails partway through.

39.14. User-Defined Operators

Every operator is “syntactic sugar” for a call to an underlying function that does the real work; so you must first create the underlying function before you can create the operator. However, an operator is *not merely* syntactic sugar, because it carries additional information that helps the query planner optimize queries that use the operator. The next section will be devoted to explaining that additional information.

Postgres Pro supports prefix and infix operators. Operators can be overloaded; that is, the same operator name can be used for different operators that have different numbers and types of operands. When a query is executed, the system determines the operator to call from the number and types of the provided operands.

Here is an example of creating an operator for adding two complex numbers. We assume we've already created the definition of type `complex` (see [Section 39.13](#)). First we need a function that does the work, then we can define the operator:

```
CREATE FUNCTION complex_add(complex, complex)
    RETURNS complex
    AS 'filename', 'complex_add'
    LANGUAGE C IMMUTABLE STRICT;

CREATE OPERATOR + (
    leftarg = complex,
    rightarg = complex,
    function = complex_add,
    commutator = +
);
```

Now we could execute a query like this:

```
SELECT (a + b) AS c FROM test_complex;
```

c

```
-----
(5.2,6.05)
(133.42,144.95)
```

We've shown how to create a binary operator here. To create a prefix operator, just omit the `leftarg`. The `function` clause and the argument clauses are the only required items in `CREATE OPERATOR`. The `commutator` clause shown in the example is an optional hint to the query optimizer. Further details about `commutator` and other optimizer hints appear in the next section.

39.15. Operator Optimization Information

A Postgres Pro operator definition can include several optional clauses that tell the system useful things about how the operator behaves. These clauses should be provided whenever appropriate, because they can make for considerable speedups in execution of queries that use the operator. But if you provide them, you must be sure that they are right! Incorrect use of an optimization clause can result in slow queries, subtly wrong output, or other Bad Things. You can always leave out an optimization clause if you are not sure about it; the only consequence is that queries might run slower than they need to.

Additional optimization clauses might be added in future versions of Postgres Pro. The ones described here are all the ones that release 17.5.1 understands.

It is also possible to attach a planner support function to the function that underlies an operator, providing another way of telling the system about the behavior of the operator. See [Section 39.11](#) for more information.

39.15.1. COMMUTATOR

The `COMMUTATOR` clause, if provided, names an operator that is the commutator of the operator being defined. We say that operator A is the commutator of operator B if $(x \text{ A } y)$ equals $(y \text{ B } x)$ for all possible input values x, y . Notice that B is also the commutator of A. For example, operators `<` and `>` for a particular data type are usually each others' commutators, and operator `+` is usually commutative with itself. But operator `-` is usually not commutative with anything.

The left operand type of a commutable operator is the same as the right operand type of its commutator, and vice versa. So the name of the commutator operator is all that Postgres Pro needs to be given to look up the commutator, and that's all that needs to be provided in the `COMMUTATOR` clause.

It's critical to provide commutator information for operators that will be used in indexes and join clauses, because this allows the query optimizer to “flip around” such a clause to the forms needed for different plan types. For example, consider a query with a WHERE clause like `tab1.x = tab2.y`, where `tab1.x` and `tab2.y` are of a user-defined type, and suppose that `tab2.y` is indexed. The optimizer cannot generate an index scan unless it can determine how to flip the clause around to `tab2.y = tab1.x`, because the index-scan machinery expects to see the indexed column on the left of the operator it is given. Postgres Pro will *not* simply assume that this is a valid transformation — the creator of the `=` operator must specify that it is valid, by marking the operator with commutator information.

39.15.2. NEGATOR

The `NEGATOR` clause, if provided, names an operator that is the negator of the operator being defined. We say that operator A is the negator of operator B if both return Boolean results and $(x \text{ A } y)$ equals `NOT (x B y)` for all possible inputs x, y . Notice that B is also the negator of A. For example, `<` and `>=` are a negator pair for most data types. An operator can never validly be its own negator.

Unlike commutators, a pair of unary operators could validly be marked as each other's negators; that would mean $(A \text{ x})$ equals `NOT (B x)` for all x .

An operator's negator must have the same left and/or right operand types as the operator to be defined, so just as with `COMMUTATOR`, only the operator name need be given in the `NEGATOR` clause.

Providing a negator is very helpful to the query optimizer since it allows expressions like `NOT (x = y)` to be simplified into `x <> y`. This comes up more often than you might think, because `NOT` operations can be inserted as a consequence of other rearrangements.

39.15.3. RESTRICT

The `RESTRICT` clause, if provided, names a restriction selectivity estimation function for the operator. (Note that this is a function name, not an operator name.) `RESTRICT` clauses only make sense for binary operators that return `boolean`. The idea behind a restriction selectivity estimator is to guess what fraction of the rows in a table will satisfy a `WHERE`-clause condition of the form:

```
column OP constant
```

for the current operator and a particular constant value. This assists the optimizer by giving it some idea of how many rows will be eliminated by `WHERE` clauses that have this form. (What happens if the constant is on the left, you might be wondering? Well, that's one of the things that `COMMUTATOR` is for...)

Writing new restriction selectivity estimation functions is far beyond the scope of this chapter, but fortunately you can usually just use one of the system's standard estimators for many of your own operators. These are the standard restriction estimators:

```
eqsel for =
neqsel for <>
scalarltsel for <
scalarlesel for <=
scalargtsel for >
scalargesel for >=
```

You can frequently get away with using either `eqsel` or `neqsel` for operators that have very high or very low selectivity, even if they aren't really equality or inequality. For example, the approximate-equality geometric operators use `eqsel` on the assumption that they'll usually only match a small fraction of the entries in a table.

You can use `scalarltsel`, `scalarlesel`, `scalargtsel` and `scalargesel` for comparisons on data types that have some sensible means of being converted into numeric scalars for range comparisons.

Another useful built-in selectivity estimation function is `matchingsel`, which will work for almost any binary operator, if standard MCV and/or histogram statistics are collected for the input data type(s). Its default estimate is set to twice the default estimate used in `eqsel`, making it most suitable for comparison operators that are somewhat less strict than equality. (Or you could call the underlying `generic_restriction_selectivity` function, providing a different default estimate.)

39.15.4. JOIN

The `JOIN` clause, if provided, names a join selectivity estimation function for the operator. (Note that this is a function name, not an operator name.) `JOIN` clauses only make sense for binary operators that return `boolean`. The idea behind a join selectivity estimator is to guess what fraction of the rows in a pair of tables will satisfy a `WHERE`-clause condition of the form:

```
table1.column1 OP table2.column2
```

for the current operator. As with the `RESTRICT` clause, this helps the optimizer very substantially by letting it figure out which of several possible join sequences is likely to take the least work.

As before, this chapter will make no attempt to explain how to write a join selectivity estimator function, but will just suggest that you use one of the standard estimators if one is applicable:

```
eqjoinssel for =
neqjoinssel for <>
scalarltjoinssel for <
scalarlejoinssel for <=
scalargtjoinssel for >
```

`scalargejoin` for `>=`
`matchingjoin` for generic matching operators
`areajoin` for 2D area-based comparisons
`positionjoin` for 2D position-based comparisons
`containjoin` for 2D containment-based comparisons

39.15.5. HASHES

The `HASHES` clause, if present, tells the system that it is permissible to use the hash join method for a join based on this operator. `HASHES` only makes sense for a binary operator that returns `boolean`, and in practice the operator must represent equality for some data type or pair of data types.

The assumption underlying hash join is that the join operator can only return true for pairs of left and right values that hash to the same hash code. If two values get put in different hash buckets, the join will never compare them at all, implicitly assuming that the result of the join operator must be false. So it never makes sense to specify `HASHES` for operators that do not represent some form of equality. In most cases it is only practical to support hashing for operators that take the same data type on both sides. However, sometimes it is possible to design compatible hash functions for two or more data types; that is, functions that will generate the same hash codes for “equal” values, even though the values have different representations. For example, it's fairly simple to arrange this property when hashing integers of different widths.

To be marked `HASHES`, the join operator must appear in a hash index operator family. This is not enforced when you create the operator, since of course the referencing operator family couldn't exist yet. But attempts to use the operator in hash joins will fail at run time if no such operator family exists. The system needs the operator family to find the data-type-specific hash function(s) for the operator's input data type(s). Of course, you must also create suitable hash functions before you can create the operator family.

Care should be exercised when preparing a hash function, because there are machine-dependent ways in which it might fail to do the right thing. For example, if your data type is a structure in which there might be uninteresting pad bits, you cannot simply pass the whole structure to `hash_any`. (Unless you write your other operators and functions to ensure that the unused bits are always zero, which is the recommended strategy.) Another example is that on machines that meet the IEEE floating-point standard, negative zero and positive zero are different values (different bit patterns) but they are defined to compare equal. If a float value might contain negative zero then extra steps are needed to ensure it generates the same hash value as positive zero.

A hash-joinable operator must have a commutator (itself if the two operand data types are the same, or a related equality operator if they are different) that appears in the same operator family. If this is not the case, planner errors might occur when the operator is used. Also, it is a good idea (but not strictly required) for a hash operator family that supports multiple data types to provide equality operators for every combination of the data types; this allows better optimization.

Note

The function underlying a hash-joinable operator must be marked immutable or stable. If it is volatile, the system will never attempt to use the operator for a hash join.

Note

If a hash-joinable operator has an underlying function that is marked strict, the function must also be complete: that is, it should return true or false, never null, for any two nonnull inputs. If this rule is not followed, hash-optimization of `IN` operations might generate wrong results. (Specifically, `IN` might return false where the correct answer according to the standard would be null; or it might yield an error complaining that it wasn't prepared for a null result.)

39.15.6. MERGES

The `MERGES` clause, if present, tells the system that it is permissible to use the merge-join method for a join based on this operator. `MERGES` only makes sense for a binary operator that returns `boolean`, and in practice the operator must represent equality for some data type or pair of data types.

Merge join is based on the idea of sorting the left- and right-hand tables into order and then scanning them in parallel. So, both data types must be capable of being fully ordered, and the join operator must be one that can only succeed for pairs of values that fall at the “same place” in the sort order. In practice this means that the join operator must behave like equality. But it is possible to merge-join two distinct data types so long as they are logically compatible. For example, the `smallint-versus-integer` equality operator is merge-joinable. We only need sorting operators that will bring both data types into a logically compatible sequence.

To be marked `MERGES`, the join operator must appear as an equality member of a `btree` index operator family. This is not enforced when you create the operator, since of course the referencing operator family couldn't exist yet. But the operator will not actually be used for merge joins unless a matching operator family can be found. The `MERGES` flag thus acts as a hint to the planner that it's worth looking for a matching operator family.

A merge-joinable operator must have a commutator (itself if the two operand data types are the same, or a related equality operator if they are different) that appears in the same operator family. If this is not the case, planner errors might occur when the operator is used. Also, it is a good idea (but not strictly required) for a `btree` operator family that supports multiple data types to provide equality operators for every combination of the data types; this allows better optimization.

Note

The function underlying a merge-joinable operator must be marked immutable or stable. If it is volatile, the system will never attempt to use the operator for a merge join.

39.16. Interfacing Extensions to Indexes

The procedures described thus far let you define new types, new functions, and new operators. However, we cannot yet define an index on a column of a new data type. To do this, we must define an *operator class* for the new data type. Later in this section, we will illustrate this concept in an example: a new operator class for the B-tree index method that stores and sorts complex numbers in ascending absolute value order.

Operator classes can be grouped into *operator families* to show the relationships between semantically compatible classes. When only a single data type is involved, an operator class is sufficient, so we'll focus on that case first and then return to operator families.

39.16.1. Index Methods and Operator Classes

Operator classes are associated with an index access method, such as `B-Tree` or `GIN`. Custom index access method may be defined with `CREATE ACCESS METHOD`. See [Chapter 63](#) for details.

The routines for an index method do not directly know anything about the data types that the index method will operate on. Instead, an *operator class* identifies the set of operations that the index method needs to use to work with a particular data type. Operator classes are so called because one thing they specify is the set of `WHERE`-clause operators that can be used with an index (i.e., can be converted into an index-scan qualification). An operator class can also specify some *support function* that are needed by the internal operations of the index method, but do not directly correspond to any `WHERE`-clause operator that can be used with the index.

It is possible to define multiple operator classes for the same data type and index method. By doing this, multiple sets of indexing semantics can be defined for a single data type. For example, a B-tree index requires a sort ordering to be defined for each data type it works on. It might be useful for a

complex-number data type to have one B-tree operator class that sorts the data by complex absolute value, another that sorts by real part, and so on. Typically, one of the operator classes will be deemed most commonly useful and will be marked as the default operator class for that data type and index method.

The same operator class name can be used for several different index methods (for example, both B-tree and hash index methods have operator classes named `int4_ops`), but each such class is an independent entity and must be defined separately.

39.16.2. Index Method Strategies

The operators associated with an operator class are identified by “strategy numbers”, which serve to identify the semantics of each operator within the context of its operator class. For example, B-trees impose a strict ordering on keys, lesser to greater, and so operators like “less than” and “greater than or equal to” are interesting with respect to a B-tree. Because Postgres Pro allows the user to define operators, Postgres Pro cannot look at the name of an operator (e.g., `<` or `>=`) and tell what kind of comparison it is. Instead, the index method defines a set of “strategies”, which can be thought of as generalized operators. Each operator class specifies which actual operator corresponds to each strategy for a particular data type and interpretation of the index semantics.

The B-tree index method defines six strategies, shown in [Table 39.3](#).

Table 39.3. B-Tree Strategies

Operation	Strategy Number
less than	1
less than or equal	2
equal	3
greater than or equal	4
greater than	5
distance	6

Hash indexes support only equality comparisons, and so they use only one strategy, shown in [Table 39.4](#).

Table 39.4. Hash Strategies

Operation	Strategy Number
equal	1

GiST indexes are more flexible: they do not have a fixed set of strategies at all. Instead, the “consistency” support routine of each particular GiST operator class interprets the strategy numbers however it likes. As an example, several of the built-in GiST index operator classes index two-dimensional geometric objects, providing the “R-tree” strategies shown in [Table 39.5](#). Four of these are true two-dimensional tests (overlaps, same, contains, contained by); four of them consider only the X direction; and the other four provide the same tests in the Y direction.

Table 39.5. GiST Two-Dimensional “R-tree” Strategies

Operation	Strategy Number
strictly left of	1
does not extend to right of	2
overlaps	3
does not extend to left of	4
strictly right of	5
same	6
contains	7

Operation	Strategy Number
contained by	8
does not extend above	9
strictly below	10
strictly above	11
does not extend below	12

SP-GiST indexes are similar to GiST indexes in flexibility: they don't have a fixed set of strategies. Instead the support routines of each operator class interpret the strategy numbers according to the operator class's definition. As an example, the strategy numbers used by the built-in operator classes for points are shown in [Table 39.6](#).

Table 39.6. SP-GiST Point Strategies

Operation	Strategy Number
strictly left of	1
strictly right of	5
same	6
contained by	8
strictly below	10
strictly above	11

GIN indexes are similar to GiST and SP-GiST indexes, in that they don't have a fixed set of strategies either. Instead the support routines of each operator class interpret the strategy numbers according to the operator class's definition. As an example, the strategy numbers used by the built-in operator class for arrays are shown in [Table 39.7](#).

Table 39.7. GIN Array Strategies

Operation	Strategy Number
overlap	1
contains	2
is contained by	3
equal	4

BRIN indexes are similar to GiST, SP-GiST and GIN indexes in that they don't have a fixed set of strategies either. Instead the support routines of each operator class interpret the strategy numbers according to the operator class's definition. As an example, the strategy numbers used by the built-in `Minmax` operator classes are shown in [Table 39.8](#).

Table 39.8. BRIN Minmax Strategies

Operation	Strategy Number
less than	1
less than or equal	2
equal	3
greater than or equal	4
greater than	5

Notice that all the operators listed above return Boolean values. In practice, all operators defined as index method search operators must return type `boolean`, since they must appear at the top level of a `WHERE` clause to be used with an index. (Some index access methods also support *ordering operators*, which typically don't return Boolean values; that feature is discussed in [Section 39.16.7](#).)

39.16.3. Index Method Support Routines

Strategies aren't usually enough information for the system to figure out how to use an index. In practice, the index methods require additional support routines in order to work. For example, the B-tree index method must be able to compare two keys and determine whether one is greater than, equal to, or less than the other. Similarly, the hash index method must be able to compute hash codes for key values. These operations do not correspond to operators used in qualifications in SQL commands; they are administrative routines used by the index methods, internally.

Just as with strategies, the operator class identifies which specific functions should play each of these roles for a given data type and semantic interpretation. The index method defines the set of functions it needs, and the operator class identifies the correct functions to use by assigning them to the “support function numbers” specified by the index method.

Additionally, some opclasses allow users to specify parameters which control their behavior. Each builtin index access method has an optional `options` support function, which defines a set of opclass-specific parameters.

B-trees require a comparison support function, and allow four additional support functions to be supplied at the operator class author's option, as shown in [Table 39.9](#). The requirements for these support functions are explained further in [Section 65.1.3](#).

Table 39.9. B-Tree Support Functions

Function	Support Number
Compare two keys and return an integer less than zero, zero, or greater than zero, indicating whether the first key is less than, equal to, or greater than the second	1
Return the addresses of C-callable sort support function(s) (optional)	2
Compare a test value to a base value plus/minus an offset, and return true or false according to the comparison result (optional)	3
Determine if it is safe for indexes that use the operator class to apply the btree deduplication optimization (optional)	4
Define options that are specific to this operator class (optional)	5

Hash indexes require one support function, and allow two additional ones to be supplied at the operator class author's option, as shown in [Table 39.10](#).

Table 39.10. Hash Support Functions

Function	Support Number
Compute the 32-bit hash value for a key	1
Compute the 64-bit hash value for a key given a 64-bit salt; if the salt is 0, the low 32 bits of the result must match the value that would have been computed by function 1 (optional)	2
Define options that are specific to this operator class (optional)	3

GiST indexes have eleven support functions, six of which are optional, as shown in [Table 39.11](#). (For more information see [Section 65.2](#).)

Table 39.11. GiST Support Functions

Function	Description	Support Number
<code>consistent</code>	determine whether key satisfies the query qualifier	1

Function	Description	Support Number
union	compute union of a set of keys	2
compress	compute a compressed representation of a key or value to be indexed (optional)	3
decompress	compute a decompressed representation of a compressed key (optional)	4
penalty	compute penalty for inserting new key into subtree with given subtree's key	5
picksplit	determine which entries of a page are to be moved to the new page and compute the union keys for resulting pages	6
same	compare two keys and return true if they are equal	7
distance	determine distance from key to query value (optional)	8
fetch	compute original representation of a compressed key for index-only scans (optional)	9
options	define options that are specific to this operator class (optional)	10
sortsupport	provide a sort comparator to be used in fast index builds (optional)	11

SP-GiST indexes have six support functions, one of which is optional, as shown in [Table 39.12](#). (For more information see [Section 65.3](#).)

Table 39.12. SP-GiST Support Functions

Function	Description	Support Number
config	provide basic information about the operator class	1
choose	determine how to insert a new value into an inner tuple	2
picksplit	determine how to partition a set of values	3
inner_consistent	determine which sub-partitions need to be searched for a query	4
leaf_consistent	determine whether key satisfies the query qualifier	5
options	define options that are specific to this operator class (optional)	6

GIN indexes have seven support functions, four of which are optional, as shown in [Table 39.13](#). (For more information see [Section 65.4](#).)

Table 39.13. GIN Support Functions

Function	Description	Support Number
compare	compare two keys and return an integer less than zero, zero, or greater than zero, indicating whether the first key is less than, equal to, or greater than the second	1

Function	Description	Support Number
<code>extractValue</code>	extract keys from a value to be indexed	2
<code>extractQuery</code>	extract keys from a query condition	3
<code>consistent</code>	determine whether value matches query condition (Boolean variant) (optional if support function 6 is present)	4
<code>comparePartial</code>	compare partial key from query and key from index, and return an integer less than zero, zero, or greater than zero, indicating whether GIN should ignore this index entry, treat the entry as a match, or stop the index scan (optional)	5
<code>triConsistent</code>	determine whether value matches query condition (ternary variant) (optional if support function 4 is present)	6
<code>options</code>	define options that are specific to this operator class (optional)	7

BRIN indexes have five basic support functions, one of which is optional, as shown in [Table 39.14](#). Some versions of the basic functions require additional support functions to be provided. (For more information see [Section 65.5.3](#).)

Table 39.14. BRIN Support Functions

Function	Description	Support Number
<code>opcInfo</code>	return internal information describing the indexed columns' summary data	1
<code>add_value</code>	add a new value to an existing summary index tuple	2
<code>consistent</code>	determine whether value matches query condition	3
<code>union</code>	compute union of two summary tuples	4
<code>options</code>	define options that are specific to this operator class (optional)	5

Unlike search operators, support functions return whichever data type the particular index method expects; for example in the case of the comparison function for B-trees, a signed integer. The number and types of the arguments to each support function are likewise dependent on the index method. For B-tree and hash the comparison and hashing support functions take the same input data types as do the operators included in the operator class, but this is not the case for most GiST, SP-GiST, GIN, and BRIN support functions.

39.16.4. An Example

Now that we have seen the ideas, here is the promised example of creating a new operator class. The operator class encapsulates operators that sort complex numbers in absolute value order, so we choose the name `complex_abs_ops`. First, we need a set of operators. The procedure for defining operators was discussed in [Section 39.14](#). For an operator class on B-trees, the operators we require are:

- absolute-value less-than (strategy 1)
- absolute-value less-than-or-equal (strategy 2)
- absolute-value equal (strategy 3)
- absolute-value greater-than-or-equal (strategy 4)
- absolute-value greater-than (strategy 5)

The least error-prone way to define a related set of comparison operators is to write the B-tree comparison support function first, and then write the other functions as one-line wrappers around the support function. This reduces the odds of getting inconsistent results for corner cases. Following this approach, we first write:

```
#define Mag(c)    ((c)->x*(c)->x + (c)->y*(c)->y)

static int
complex_abs_cmp_internal(Complex *a, Complex *b)
{
    double      amag = Mag(a),
               bmag = Mag(b);

    if (amag < bmag)
        return -1;
    if (amag > bmag)
        return 1;
    return 0;
}
```

Now the less-than function looks like:

```
PG_FUNCTION_INFO_V1(complex_abs_lt);

Datum
complex_abs_lt(PG_FUNCTION_ARGS)
{
    Complex      *a = (Complex *) PG_GETARG_POINTER(0);
    Complex      *b = (Complex *) PG_GETARG_POINTER(1);

    PG_RETURN_BOOL(complex_abs_cmp_internal(a, b) < 0);
}
```

The other four functions differ only in how they compare the internal function's result to zero.

Next we declare the functions and the operators based on the functions to SQL:

```
CREATE FUNCTION complex_abs_lt(complex, complex) RETURNS bool
    AS 'filename', 'complex_abs_lt'
    LANGUAGE C IMMUTABLE STRICT;

CREATE OPERATOR < (
    leftarg = complex, rightarg = complex, procedure = complex_abs_lt,
    commutator = > , negator = >= ,
    restrict = scalarlttsel, join = scalarltjoinsel
);
```

It is important to specify the correct commutator and negator operators, as well as suitable restriction and join selectivity functions, otherwise the optimizer will be unable to make effective use of the index.

Other things worth noting are happening here:

- There can only be one operator named, say, = and taking type `complex` for both operands. In this case we don't have any other operator = for `complex`, but if we were building a practical data type we'd probably want = to be the ordinary equality operation for complex numbers (and not the equality of the absolute values). In that case, we'd need to use some other operator name for `complex_abs_eq`.
- Although Postgres Pro can cope with functions having the same SQL name as long as they have different argument data types, C can only cope with one global function having a given name. So we

shouldn't name the C function something simple like `abs_eq`. Usually it's a good practice to include the data type name in the C function name, so as not to conflict with functions for other data types.

- We could have made the SQL name of the function `abs_eq`, relying on Postgres Pro to distinguish it by argument data types from any other SQL function of the same name. To keep the example simple, we make the function have the same names at the C level and SQL level.

The next step is the registration of the support routine required by B-trees. The example C code that implements this is in the same file that contains the operator functions. This is how we declare the function:

```
CREATE FUNCTION complex_abs_cmp(complex, complex)
    RETURNS integer
    AS 'filename'
    LANGUAGE C IMMUTABLE STRICT;
```

Now that we have the required operators and support routine, we can finally create the operator class:

```
CREATE OPERATOR CLASS complex_abs_ops
    DEFAULT FOR TYPE complex USING btree AS
        OPERATOR          1          < ,
        OPERATOR          2          <= ,
        OPERATOR          3          = ,
        OPERATOR          4          >= ,
        OPERATOR          5          > ,
        FUNCTION           1          complex_abs_cmp(complex, complex);
```

And we're done! It should now be possible to create and use B-tree indexes on `complex` columns.

We could have written the operator entries more verbosely, as in:

```
OPERATOR          1          < (complex, complex) ,
```

but there is no need to do so when the operators take the same data type we are defining the operator class for.

The above example assumes that you want to make this new operator class the default B-tree operator class for the `complex` data type. If you don't, just leave out the word `DEFAULT`.

39.16.5. Operator Classes and Operator Families

So far we have implicitly assumed that an operator class deals with only one data type. While there certainly can be only one data type in a particular index column, it is often useful to index operations that compare an indexed column to a value of a different data type. Also, if there is use for a cross-data-type operator in connection with an operator class, it is often the case that the other data type has a related operator class of its own. It is helpful to make the connections between related classes explicit, because this can aid the planner in optimizing SQL queries (particularly for B-tree operator classes, since the planner contains a great deal of knowledge about how to work with them).

To handle these needs, Postgres Pro uses the concept of an *operator family*. An operator family contains one or more operator classes, and can also contain indexable operators and corresponding support functions that belong to the family as a whole but not to any single class within the family. We say that such operators and functions are “loose” within the family, as opposed to being bound into a specific class. Typically each operator class contains single-data-type operators while cross-data-type operators are loose in the family.

All the operators and functions in an operator family must have compatible semantics, where the compatibility requirements are set by the index method. You might therefore wonder why bother to single out particular subsets of the family as operator classes; and indeed for many purposes the class divisions are irrelevant and the family is the only interesting grouping. The reason for defining operator classes is that they specify how much of the family is needed to support any particular index. If there is an index

using an operator class, then that operator class cannot be dropped without dropping the index — but other parts of the operator family, namely other operator classes and loose operators, could be dropped. Thus, an operator class should be specified to contain the minimum set of operators and functions that are reasonably needed to work with an index on a specific data type, and then related but non-essential operators can be added as loose members of the operator family.

As an example, Postgres Pro has a built-in B-tree operator family `integer_ops`, which includes operator classes `int8_ops`, `int4_ops`, and `int2_ops` for indexes on `bigint` (`int8`), `integer` (`int4`), and `smallint` (`int2`) columns respectively. The family also contains cross-data-type comparison operators allowing any two of these types to be compared, so that an index on one of these types can be searched using a comparison value of another type. The family could be duplicated by these definitions:

```
CREATE OPERATOR FAMILY integer_ops USING btree;

CREATE OPERATOR CLASS int8_ops
DEFAULT FOR TYPE int8 USING btree FAMILY integer_ops AS
    -- standard int8 comparisons
    OPERATOR 1 < ,
    OPERATOR 2 <= ,
    OPERATOR 3 = ,
    OPERATOR 4 >= ,
    OPERATOR 5 > ,
    FUNCTION 1 btint8cmp(int8, int8) ,
    FUNCTION 2 btint8sortsupport(internal) ,
    FUNCTION 3 in_range(int8, int8, int8, boolean, boolean) ,
    FUNCTION 4 btequalimage(oid) ;

CREATE OPERATOR CLASS int4_ops
DEFAULT FOR TYPE int4 USING btree FAMILY integer_ops AS
    -- standard int4 comparisons
    OPERATOR 1 < ,
    OPERATOR 2 <= ,
    OPERATOR 3 = ,
    OPERATOR 4 >= ,
    OPERATOR 5 > ,
    FUNCTION 1 btint4cmp(int4, int4) ,
    FUNCTION 2 btint4sortsupport(internal) ,
    FUNCTION 3 in_range(int4, int4, int4, boolean, boolean) ,
    FUNCTION 4 btequalimage(oid) ;

CREATE OPERATOR CLASS int2_ops
DEFAULT FOR TYPE int2 USING btree FAMILY integer_ops AS
    -- standard int2 comparisons
    OPERATOR 1 < ,
    OPERATOR 2 <= ,
    OPERATOR 3 = ,
    OPERATOR 4 >= ,
    OPERATOR 5 > ,
    FUNCTION 1 btint2cmp(int2, int2) ,
    FUNCTION 2 btint2sortsupport(internal) ,
    FUNCTION 3 in_range(int2, int2, int2, boolean, boolean) ,
    FUNCTION 4 btequalimage(oid) ;

ALTER OPERATOR FAMILY integer_ops USING btree ADD
    -- cross-type comparisons int8 vs int2
    OPERATOR 1 < (int8, int2) ,
    OPERATOR 2 <= (int8, int2) ,
    OPERATOR 3 = (int8, int2) ,
    OPERATOR 4 >= (int8, int2) ,
```

```

OPERATOR 5 > (int8, int2) ,
FUNCTION 1 btint82cmp(int8, int2) ,

-- cross-type comparisons int8 vs int4
OPERATOR 1 < (int8, int4) ,
OPERATOR 2 <= (int8, int4) ,
OPERATOR 3 = (int8, int4) ,
OPERATOR 4 >= (int8, int4) ,
OPERATOR 5 > (int8, int4) ,
FUNCTION 1 btint84cmp(int8, int4) ,

-- cross-type comparisons int4 vs int2
OPERATOR 1 < (int4, int2) ,
OPERATOR 2 <= (int4, int2) ,
OPERATOR 3 = (int4, int2) ,
OPERATOR 4 >= (int4, int2) ,
OPERATOR 5 > (int4, int2) ,
FUNCTION 1 btint42cmp(int4, int2) ,

-- cross-type comparisons int4 vs int8
OPERATOR 1 < (int4, int8) ,
OPERATOR 2 <= (int4, int8) ,
OPERATOR 3 = (int4, int8) ,
OPERATOR 4 >= (int4, int8) ,
OPERATOR 5 > (int4, int8) ,
FUNCTION 1 btint48cmp(int4, int8) ,

-- cross-type comparisons int2 vs int8
OPERATOR 1 < (int2, int8) ,
OPERATOR 2 <= (int2, int8) ,
OPERATOR 3 = (int2, int8) ,
OPERATOR 4 >= (int2, int8) ,
OPERATOR 5 > (int2, int8) ,
FUNCTION 1 btint28cmp(int2, int8) ,

-- cross-type comparisons int2 vs int4
OPERATOR 1 < (int2, int4) ,
OPERATOR 2 <= (int2, int4) ,
OPERATOR 3 = (int2, int4) ,
OPERATOR 4 >= (int2, int4) ,
OPERATOR 5 > (int2, int4) ,
FUNCTION 1 btint24cmp(int2, int4) ,

-- cross-type in_range functions
FUNCTION 3 in_range(int4, int4, int8, boolean, boolean) ,
FUNCTION 3 in_range(int4, int4, int2, boolean, boolean) ,
FUNCTION 3 in_range(int2, int2, int8, boolean, boolean) ,
FUNCTION 3 in_range(int2, int2, int4, boolean, boolean) ;

```

Notice that this definition “overloads” the operator strategy and support function numbers: each number occurs multiple times within the family. This is allowed so long as each instance of a particular number has distinct input data types. The instances that have both input types equal to an operator class's input type are the primary operators and support functions for that operator class, and in most cases should be declared as part of the operator class rather than as loose members of the family.

In a B-tree operator family, all the operators in the family must sort compatibly, as is specified in detail in [Section 65.1.2](#). For each operator in the family there must be a support function having the same two input data types as the operator. It is recommended that a family be complete, i.e., for each combination

of data types, all operators are included. Each operator class should include just the non-cross-type operators and support function for its data type.

To build a multiple-data-type hash operator family, compatible hash support functions must be created for each data type supported by the family. Here compatibility means that the functions are guaranteed to return the same hash code for any two values that are considered equal by the family's equality operators, even when the values are of different types. This is usually difficult to accomplish when the types have different physical representations, but it can be done in some cases. Furthermore, casting a value from one data type represented in the operator family to another data type also represented in the operator family via an implicit or binary coercion cast must not change the computed hash value. Notice that there is only one support function per data type, not one per equality operator. It is recommended that a family be complete, i.e., provide an equality operator for each combination of data types. Each operator class should include just the non-cross-type equality operator and the support function for its data type.

GiST, SP-GiST, and GIN indexes do not have any explicit notion of cross-data-type operations. The set of operators supported is just whatever the primary support functions for a given operator class can handle.

In BRIN, the requirements depends on the framework that provides the operator classes. For operator classes based on `minmax`, the behavior required is the same as for B-tree operator families: all the operators in the family must sort compatibly, and casts must not change the associated sort ordering.

Note

Prior to PostgreSQL 8.3, there was no concept of operator families, and so any cross-data-type operators intended to be used with an index had to be bound directly into the index's operator class. While this approach still works, it is deprecated because it makes an index's dependencies too broad, and because the planner can handle cross-data-type comparisons more effectively when both data types have operators in the same operator family.

39.16.6. System Dependencies on Operator Classes

Postgres Pro uses operator classes to infer the properties of operators in more ways than just whether they can be used with indexes. Therefore, you might want to create operator classes even if you have no intention of indexing any columns of your data type.

In particular, there are SQL features such as `ORDER BY` and `DISTINCT` that require comparison and sorting of values. To implement these features on a user-defined data type, Postgres Pro looks for the default B-tree operator class for the data type. The “equals” member of this operator class defines the system's notion of equality of values for `GROUP BY` and `DISTINCT`, and the sort ordering imposed by the operator class defines the default `ORDER BY` ordering.

If there is no default B-tree operator class for a data type, the system will look for a default hash operator class. But since that kind of operator class only provides equality, it is only able to support grouping not sorting.

When there is no default operator class for a data type, you will get errors like “could not identify an ordering operator” if you try to use these SQL features with the data type.

Note

In PostgreSQL versions before 7.4, sorting and grouping operations would implicitly use operators named `=`, `<`, and `>`. The new behavior of relying on default operator classes avoids having to make any assumption about the behavior of operators with particular names.

Sorting by a non-default B-tree operator class is possible by specifying the class's less-than operator in a `USING` option, for example

```
SELECT * FROM mytable ORDER BY somecol USING ~<~;
```

Alternatively, specifying the class's greater-than operator in `USING` selects a descending-order sort.

Comparison of arrays of a user-defined type also relies on the semantics defined by the type's default B-tree operator class. If there is no default B-tree operator class, but there is a default hash operator class, then array equality is supported, but not ordering comparisons.

Another SQL feature that requires even more data-type-specific knowledge is the `RANGE offset PRECEDING/FOLLOWING` framing option for window functions (see [Section 4.2.8](#)). For a query such as

```
SELECT sum(x) OVER (ORDER BY x RANGE BETWEEN 5 PRECEDING AND 10 FOLLOWING)
FROM mytable;
```

it is not sufficient to know how to order by `x`; the database must also understand how to “subtract 5” or “add 10” to the current row's value of `x` to identify the bounds of the current window frame. Comparing the resulting bounds to other rows' values of `x` is possible using the comparison operators provided by the B-tree operator class that defines the `ORDER BY` ordering — but addition and subtraction operators are not part of the operator class, so which ones should be used? Hard-wiring that choice would be undesirable, because different sort orders (different B-tree operator classes) might need different behavior. Therefore, a B-tree operator class can specify an *in_range* support function that encapsulates the addition and subtraction behaviors that make sense for its sort order. It can even provide more than one *in_range* support function, in case there is more than one data type that makes sense to use as the offset in `RANGE` clauses. If the B-tree operator class associated with the window's `ORDER BY` clause does not have a matching *in_range* support function, the `RANGE offset PRECEDING/FOLLOWING` option is not supported.

Another important point is that an equality operator that appears in a hash operator family is a candidate for hash joins, hash aggregation, and related optimizations. The hash operator family is essential here since it identifies the hash function(s) to use.

39.16.7. Ordering Operators

Some index access methods (currently, only B-tree, GiST and SP-GiST) support the concept of *ordering operators*. What we have been discussing so far are *search operators*. A search operator is one for which the index can be searched to find all rows satisfying `WHERE indexed_column operator constant`. Note that nothing is promised about the order in which the matching rows will be returned. In contrast, an ordering operator does not restrict the set of rows that can be returned, but instead determines their order. An ordering operator is one for which the index can be scanned to return rows in the order represented by `ORDER BY indexed_column operator constant`. The reason for defining ordering operators that way is that it supports nearest-neighbor searches, if the operator is one that measures distance. For example, a query like

```
SELECT * FROM places ORDER BY location <-> point '(101,456)' LIMIT 10;
```

finds the ten places closest to a given target point. A GiST index on the location column can do this efficiently because `<->` is an ordering operator.

While search operators have to return Boolean results, ordering operators usually return some other type, such as float or numeric for distances. This type is normally not the same as the data type being indexed. To avoid hard-wiring assumptions about the behavior of different data types, the definition of an ordering operator is required to name a B-tree operator family that specifies the sort ordering of the result data type. As was stated in the previous section, B-tree operator families define Postgres Pro's notion of ordering, so this is a natural representation. Since the point `<->` operator returns `float8`, it could be specified in an operator class creation command like this:

```
OPERATOR 15      <-> (point, point) FOR ORDER BY float_ops
```

where `float_ops` is the built-in operator family that includes operations on `float8`. This declaration states that the index is able to return rows in order of increasing values of the `<->` operator.

39.16.8. Special Features of Operator Classes

There are two special features of operator classes that we have not discussed yet, mainly because they are not useful with the most commonly used index methods.

Normally, declaring an operator as a member of an operator class (or family) means that the index method can retrieve exactly the set of rows that satisfy a `WHERE` condition using the operator. For example:

```
SELECT * FROM table WHERE integer_column < 4;
```

can be satisfied exactly by a B-tree index on the integer column. But there are cases where an index is useful as an inexact guide to the matching rows. For example, if a GiST index stores only bounding boxes for geometric objects, then it cannot exactly satisfy a `WHERE` condition that tests overlap between nonrectangular objects such as polygons. Yet we could use the index to find objects whose bounding box overlaps the bounding box of the target object, and then do the exact overlap test only on the objects found by the index. If this scenario applies, the index is said to be “lossy” for the operator. Lossy index searches are implemented by having the index method return a *recheck* flag when a row might or might not really satisfy the query condition. The core system will then test the original query condition on the retrieved row to see whether it should be returned as a valid match. This approach works if the index is guaranteed to return all the required rows, plus perhaps some additional rows, which can be eliminated by performing the original operator invocation. The index methods that support lossy searches (currently, GiST, SP-GiST and GIN) allow the support functions of individual operator classes to set the recheck flag, and so this is essentially an operator-class feature.

Consider again the situation where we are storing in the index only the bounding box of a complex object such as a polygon. In this case there's not much value in storing the whole polygon in the index entry — we might as well store just a simpler object of type `box`. This situation is expressed by the `STORAGE` option in `CREATE OPERATOR CLASS`: we'd write something like:

```
CREATE OPERATOR CLASS polygon_ops
    DEFAULT FOR TYPE polygon USING gist AS
    ...
    STORAGE box;
```

At present, only the GiST, SP-GiST, GIN and BRIN index methods support a `STORAGE` type that's different from the column data type. The GiST `compress` and `decompress` support routines must deal with data-type conversion when `STORAGE` is used. SP-GiST likewise requires a `compress` support function to convert to the storage type, when that is different; if an SP-GiST opclass also supports retrieving data, the reverse conversion must be handled by the `consistent` function. In GIN, the `STORAGE` type identifies the type of the “key” values, which normally is different from the type of the indexed column — for example, an operator class for integer-array columns might have keys that are just integers. The GIN `extractValue` and `extractQuery` support routines are responsible for extracting keys from indexed values. BRIN is similar to GIN: the `STORAGE` type identifies the type of the stored summary values, and operator classes' support procedures are responsible for interpreting the summary values correctly.

39.17. Packaging Related Objects into an Extension

A useful extension to Postgres Pro typically includes multiple SQL objects; for example, a new data type will require new functions, new operators, and probably new index operator classes. It is helpful to collect all these objects into a single package to simplify database management. Postgres Pro calls such a package an *extension*. To define an extension, you need at least a *script file* that contains the SQL commands to create the extension's objects, and a *control file* that specifies a few basic properties of the extension itself. If the extension includes C code, there will typically also be a shared library file into which the C code has been built. Once you have these files, a simple `CREATE EXTENSION` command loads the objects into your database.

The main advantage of using an extension, rather than just running the SQL script to load a bunch of “loose” objects into your database, is that Postgres Pro will then understand that the objects of the extension go together. You can drop all the objects with a single `DROP EXTENSION` command (no need

to maintain a separate “uninstall” script). Even more useful, `pg_dump` knows that it should not dump the individual member objects of the extension — it will just include a `CREATE EXTENSION` command in dumps, instead. This vastly simplifies migration to a new version of the extension that might contain more or different objects than the old version. Note however that you must have the extension's control, script, and other files available when loading such a dump into a new database.

Postgres Pro will not let you drop an individual object contained in an extension, except by dropping the whole extension. Also, while you can change the definition of an extension member object (for example, via `CREATE OR REPLACE FUNCTION` for a function), bear in mind that the modified definition will not be dumped by `pg_dump`. Such a change is usually only sensible if you concurrently make the same change in the extension's script file. (But there are special provisions for tables containing configuration data; see [Section 39.17.3](#).) In production situations, it's generally better to create an extension update script to perform changes to extension member objects.

The extension script may set privileges on objects that are part of the extension, using `GRANT` and `REVOKE` statements. The final set of privileges for each object (if any are set) will be stored in the `pg_init_privs` system catalog. When `pg_dump` is used, the `CREATE EXTENSION` command will be included in the dump, followed by the set of `GRANT` and `REVOKE` statements necessary to set the privileges on the objects to what they were at the time the dump was taken.

Postgres Pro does not currently support extension scripts issuing `CREATE POLICY` or `SECURITY LABEL` statements. These are expected to be set after the extension has been created. All RLS policies and security labels on extension objects will be included in dumps created by `pg_dump`.

The extension mechanism also has provisions for packaging modification scripts that adjust the definitions of the SQL objects contained in an extension. For example, if version 1.1 of an extension adds one function and changes the body of another function compared to 1.0, the extension author can provide an *update script* that makes just those two changes. The `ALTER EXTENSION UPDATE` command can then be used to apply these changes and track which version of the extension is actually installed in a given database.

The kinds of SQL objects that can be members of an extension are shown in the description of [ALTER EXTENSION](#). Notably, objects that are database-cluster-wide, such as databases, roles, and tablespaces, cannot be extension members since an extension is only known within one database. (Although an extension script is not prohibited from creating such objects, if it does so they will not be tracked as part of the extension.) Also notice that while a table can be a member of an extension, its subsidiary objects such as indexes are not directly considered members of the extension. Another important point is that schemas can belong to extensions, but not vice versa: an extension as such has an unqualified name and does not exist “within” any schema. The extension's member objects, however, will belong to schemas whenever appropriate for their object types. It may or may not be appropriate for an extension to own the schema(s) its member objects are within.

If an extension's script creates any temporary objects (such as temp tables), those objects are treated as extension members for the remainder of the current session, but are automatically dropped at session end, as any temporary object would be. This is an exception to the rule that extension member objects cannot be dropped without dropping the whole extension.

39.17.1. Extension Files

The `CREATE EXTENSION` command relies on a control file for each extension, which must be named the same as the extension with a suffix of `.control`, and must be placed in the installation's `SHAREDIR/extension` directory. There must also be at least one SQL script file, which follows the naming pattern `extension--version.sql` (for example, `foo--1.0.sql` for version 1.0 of extension `foo`). By default, the script file(s) are also placed in the `SHAREDIR/extension` directory; but the control file can specify a different directory for the script file(s).

The file format for an extension control file is the same as for the `postgresql.conf` file, namely a list of `parameter_name = value` assignments, one per line. Blank lines and comments introduced by `#` are allowed. Be sure to quote any value that is not a single word or number.

A control file can set the following parameters:

`directory (string)`

The directory containing the extension's SQL script file(s). Unless an absolute path is given, the name is relative to the installation's `SHAREDIR` directory. The default behavior is equivalent to specifying `directory = 'extension'`.

`default_version (string)`

The default version of the extension (the one that will be installed if no version is specified in `CREATE EXTENSION`). Although this can be omitted, that will result in `CREATE EXTENSION` failing if no `VERSION` option appears, so you generally don't want to do that.

`comment (string)`

A comment (any string) about the extension. The comment is applied when initially creating an extension, but not during extension updates (since that might override user-added comments). Alternatively, the extension's comment can be set by writing a `COMMENT` command in the script file.

`encoding (string)`

The character set encoding used by the script file(s). This should be specified if the script files contain any non-ASCII characters. Otherwise the files will be assumed to be in the database encoding.

`module_pathname (string)`

The value of this parameter will be substituted for each occurrence of `MODULE_PATHNAME` in the script file(s). If it is not set, no substitution is made. Typically, this is set to `$libdir/shared_library_name` and then `MODULE_PATHNAME` is used in `CREATE FUNCTION` commands for C-language functions, so that the script files do not need to hard-wire the name of the shared library.

`requires (string)`

A list of names of extensions that this extension depends on, for example `requires = 'foo, bar'`. Those extensions must be installed before this one can be installed.

`no_relocate (string)`

A list of names of extensions that this extension depends on that should be barred from changing their schemas via `ALTER EXTENSION ... SET SCHEMA`. This is needed if this extension's script references the name of a required extension's schema (using the `@extschema:name@` syntax) in a way that cannot track renames.

`superuser (boolean)`

If this parameter is `true` (which is the default), only superusers can create the extension or update it to a new version (but see also `trusted`, below). If it is set to `false`, just the privileges required to execute the commands in the installation or update script are required. This should normally be set to `true` if any of the script commands require superuser privileges. (Such commands would fail anyway, but it's more user-friendly to give the error up front.)

`trusted (boolean)`

This parameter, if set to `true` (which is not the default), allows some non-superusers to install an extension that has `superuser` set to `true`. Specifically, installation will be permitted for anyone who has `CREATE` privilege on the current database. When the user executing `CREATE EXTENSION` is not a superuser but is allowed to install by virtue of this parameter, then the installation or update script is run as the bootstrap superuser, not as the calling user. This parameter is irrelevant if `superuser` is `false`. Generally, this should not be set true for extensions that could allow access to otherwise-superuser-only abilities, such as file system access. Also, marking an extension trusted requires significant extra effort to write the extension's installation and update script(s) securely; see [Section 39.17.6](#).

`relocatable` (boolean)

An extension is *relocatable* if it is possible to move its contained objects into a different schema after initial creation of the extension. The default is `false`, i.e., the extension is not relocatable. See [Section 39.17.2](#) for more information.

`schema` (string)

This parameter can only be set for non-relocatable extensions. It forces the extension to be loaded into exactly the named schema and not any other. The `schema` parameter is consulted only when initially creating an extension, not during extension updates. See [Section 39.17.2](#) for more information.

In addition to the primary control file `extension.control`, an extension can have secondary control files named in the style `extension--version.control`. If supplied, these must be located in the script file directory. Secondary control files follow the same format as the primary control file. Any parameters set in a secondary control file override the primary control file when installing or updating to that version of the extension. However, the parameters `directory` and `default_version` cannot be set in a secondary control file.

An extension's SQL script files can contain any SQL commands, except for transaction control commands (`BEGIN`, `COMMIT`, etc.) and commands that cannot be executed inside a transaction block (such as `VACUUM`). This is because the script files are implicitly executed within a transaction block.

An extension's SQL script files can also contain lines beginning with `\echo`, which will be ignored (treated as comments) by the extension mechanism. This provision is commonly used to throw an error if the script file is fed to `psql` rather than being loaded via `CREATE EXTENSION` (see example script in [Section 39.17.7](#)). Without that, users might accidentally load the extension's contents as “loose” objects rather than as an extension, a state of affairs that's a bit tedious to recover from.

If the extension script contains the string `@extowner@`, that string is replaced with the (suitably quoted) name of the user calling `CREATE EXTENSION` or `ALTER EXTENSION`. Typically this feature is used by extensions that are marked trusted to assign ownership of selected objects to the calling user rather than the bootstrap superuser. (One should be careful about doing so, however. For example, assigning ownership of a C-language function to a non-superuser would create a privilege escalation path for that user.)

While the script files can contain any characters allowed by the specified encoding, control files should contain only plain ASCII, because there is no way for Postgres Pro to know what encoding a control file is in. In practice this is only an issue if you want to use non-ASCII characters in the extension's comment. Recommended practice in that case is to not use the control file `comment` parameter, but instead use `COMMENT ON EXTENSION` within a script file to set the comment.

39.17.2. Extension Relocatability

Users often wish to load the objects contained in an extension into a different schema than the extension's author had in mind. There are three supported levels of relocatability:

- A fully relocatable extension can be moved into another schema at any time, even after it's been loaded into a database. This is done with the `ALTER EXTENSION SET SCHEMA` command, which automatically renames all the member objects into the new schema. Normally, this is only possible if the extension contains no internal assumptions about what schema any of its objects are in. Also, the extension's objects must all be in one schema to begin with (ignoring objects that do not belong to any schema, such as procedural languages). Mark a fully relocatable extension by setting `relocatable = true` in its control file.
- An extension might be relocatable during installation but not afterwards. This is typically the case if the extension's script file needs to reference the target schema explicitly, for example in setting `search_path` properties for SQL functions. For such an extension, set `relocatable = false` in its control file, and use `@extschema@` to refer to the target schema in the script file. All occurrences of this string will be replaced by the actual target schema's name (double-quoted if necessary) before the script is executed. The user can set the target schema using the `SCHEMA` option of `CREATE EXTENSION`.

- If the extension does not support relocation at all, set `relocatable = false` in its control file, and also set `schema` to the name of the intended target schema. This will prevent use of the `SCHEMA` option of `CREATE EXTENSION`, unless it specifies the same schema named in the control file. This choice is typically necessary if the extension contains internal assumptions about its schema name that can't be replaced by uses of `@extschema@`. The `@extschema@` substitution mechanism is available in this case too, although it is of limited use since the schema name is determined by the control file.

In all cases, the script file will be executed with `search_path` initially set to point to the target schema; that is, `CREATE EXTENSION` does the equivalent of this:

```
SET LOCAL search_path TO @extschema@, pg_temp;
```

This allows the objects created by the script file to go into the target schema. The script file can change `search_path` if it wishes, but that is generally undesirable. `search_path` is restored to its previous setting upon completion of `CREATE EXTENSION`.

The target schema is determined by the `schema` parameter in the control file if that is given, otherwise by the `SCHEMA` option of `CREATE EXTENSION` if that is given, otherwise the current default object creation schema (the first one in the caller's `search_path`). When the control file `schema` parameter is used, the target schema will be created if it doesn't already exist, but in the other two cases it must already exist.

If any prerequisite extensions are listed in `requires` in the control file, their target schemas are added to the initial setting of `search_path`, following the new extension's target schema. This allows their objects to be visible to the new extension's script file.

For security, `pg_temp` is automatically appended to the end of `search_path` in all cases.

Although a non-relocatable extension can contain objects spread across multiple schemas, it is usually desirable to place all the objects meant for external use into a single schema, which is considered the extension's target schema. Such an arrangement works conveniently with the default setting of `search_path` during creation of dependent extensions.

If an extension references objects belonging to another extension, it is recommended to schema-qualify those references. To do that, write `@extschema:name@` in the extension's script file, where `name` is the name of the other extension (which must be listed in this extension's `requires` list). This string will be replaced by the name (double-quoted if necessary) of that extension's target schema. Although this notation avoids the need to make hard-wired assumptions about schema names in the extension's script file, its use may embed the other extension's schema name into the installed objects of this extension. (Typically, that happens when `@extschema:name@` is used inside a string literal, such as a function body or a `search_path` setting. In other cases, the object reference is reduced to an OID during parsing and does not require subsequent lookups.) If the other extension's schema name is so embedded, you should prevent the other extension from being relocated after yours is installed, by adding the name of the other extension to this one's `no_relocate` list.

39.17.3. Extension Configuration Tables

Some extensions include configuration tables, which contain data that might be added or changed by the user after installation of the extension. Ordinarily, if a table is part of an extension, neither the table's definition nor its content will be dumped by `pg_dump`. But that behavior is undesirable for a configuration table; any data changes made by the user need to be included in dumps, or the extension will behave differently after a dump and restore.

To solve this problem, an extension's script file can mark a table or a sequence it has created as a configuration relation, which will cause `pg_dump` to include the table's or the sequence's contents (not its definition) in dumps. To do that, call the function `pg_extension_config_dump(regclass, text)` after creating the table or the sequence, for example

```
CREATE TABLE my_config (key text, value text);
CREATE SEQUENCE my_config_seq;
```

```
SELECT pg_catalog.pg_extension_config_dump('my_config', '');
SELECT pg_catalog.pg_extension_config_dump('my_config_seq', '');
```

Any number of tables or sequences can be marked this way. Sequences associated with `serial` or `bigserial` columns can be marked as well.

When the second argument of `pg_extension_config_dump` is an empty string, the entire contents of the table are dumped by `pg_dump`. This is usually only correct if the table is initially empty as created by the extension script. If there is a mixture of initial data and user-provided data in the table, the second argument of `pg_extension_config_dump` provides a `WHERE` condition that selects the data to be dumped. For example, you might do

```
CREATE TABLE my_config (key text, value text, standard_entry boolean);

SELECT pg_catalog.pg_extension_config_dump('my_config', 'WHERE NOT standard_entry');
```

and then make sure that `standard_entry` is true only in the rows created by the extension's script.

For sequences, the second argument of `pg_extension_config_dump` has no effect.

More complicated situations, such as initially-provided rows that might be modified by users, can be handled by creating triggers on the configuration table to ensure that modified rows are marked correctly.

You can alter the filter condition associated with a configuration table by calling `pg_extension_config_dump` again. (This would typically be useful in an extension update script.) The only way to mark a table as no longer a configuration table is to dissociate it from the extension with `ALTER EXTENSION ... DROP TABLE`.

Note that foreign key relationships between these tables will dictate the order in which the tables are dumped out by `pg_dump`. Specifically, `pg_dump` will attempt to dump the referenced-by table before the referencing table. As the foreign key relationships are set up at `CREATE EXTENSION` time (prior to data being loaded into the tables) circular dependencies are not supported. When circular dependencies exist, the data will still be dumped out but the dump will not be able to be restored directly and user intervention will be required.

Sequences associated with `serial` or `bigserial` columns need to be directly marked to dump their state. Marking their parent relation is not enough for this purpose.

39.17.4. Extension Updates

One advantage of the extension mechanism is that it provides convenient ways to manage updates to the SQL commands that define an extension's objects. This is done by associating a version name or number with each released version of the extension's installation script. In addition, if you want users to be able to update their databases dynamically from one version to the next, you should provide *update scripts* that make the necessary changes to go from one version to the next. Update scripts have names following the pattern `extension--old_version--target_version.sql` (for example, `foo--1.0--1.1.sql` contains the commands to modify version 1.0 of extension `foo` into version 1.1).

Given that a suitable update script is available, the command `ALTER EXTENSION UPDATE` will update an installed extension to the specified new version. The update script is run in the same environment that `CREATE EXTENSION` provides for installation scripts: in particular, `search_path` is set up in the same way, and any new objects created by the script are automatically added to the extension. Also, if the script chooses to drop extension member objects, they are automatically dissociated from the extension.

If an extension has secondary control files, the control parameters that are used for an update script are those associated with the script's target (new) version.

`ALTER EXTENSION` is able to execute sequences of update script files to achieve a requested update. For example, if only `foo--1.0--1.1.sql` and `foo--1.1--2.0.sql` are available, `ALTER EXTENSION` will apply them in sequence if an update to version 2.0 is requested when 1.0 is currently installed.

Postgres Pro doesn't assume anything about the properties of version names: for example, it does not know whether 1.1 follows 1.0. It just matches up the available version names and follows the path that requires applying the fewest update scripts. (A version name can actually be any string that doesn't contain -- or leading or trailing -.)

Sometimes it is useful to provide “downgrade” scripts, for example `foo--1.1--1.0.sql` to allow reverting the changes associated with version 1.1. If you do that, be careful of the possibility that a downgrade script might unexpectedly get applied because it yields a shorter path. The risky case is where there is a “fast path” update script that jumps ahead several versions as well as a downgrade script to the fast path's start point. It might take fewer steps to apply the downgrade and then the fast path than to move ahead one version at a time. If the downgrade script drops any irreplaceable objects, this will yield undesirable results.

To check for unexpected update paths, use this command:

```
SELECT * FROM pg_extension_update_paths('extension_name');
```

This shows each pair of distinct known version names for the specified extension, together with the update path sequence that would be taken to get from the source version to the target version, or `NULL` if there is no available update path. The path is shown in textual form with -- separators. You can use `regexp_split_to_array(path, '--')` if you prefer an array format.

39.17.5. Installing Extensions Using Update Scripts

An extension that has been around for awhile will probably exist in several versions, for which the author will need to write update scripts. For example, if you have released a `foo` extension in versions 1.0, 1.1, and 1.2, there should be update scripts `foo--1.0--1.1.sql` and `foo--1.1--1.2.sql`. Before Postgres Pro 10, it was necessary to also create new script files `foo--1.1.sql` and `foo--1.2.sql` that directly build the newer extension versions, or else the newer versions could not be installed directly, only by installing 1.0 and then updating. That was tedious and duplicative, but now it's unnecessary, because `CREATE EXTENSION` can follow update chains automatically. For example, if only the script files `foo--1.0.sql`, `foo--1.0--1.1.sql`, and `foo--1.1--1.2.sql` are available then a request to install version 1.2 is honored by running those three scripts in sequence. The processing is the same as if you'd first installed 1.0 and then updated to 1.2. (As with `ALTER EXTENSION UPDATE`, if multiple pathways are available then the shortest is preferred.) Arranging an extension's script files in this style can reduce the amount of maintenance effort needed to produce small updates.

If you use secondary (version-specific) control files with an extension maintained in this style, keep in mind that each version needs a control file even if it has no stand-alone installation script, as that control file will determine how the implicit update to that version is performed. For example, if `foo--1.0.control` specifies `requires = 'bar'` but `foo`'s other control files do not, the extension's dependency on `bar` will be dropped when updating from 1.0 to another version.

39.17.6. Security Considerations for Extensions

Widely-distributed extensions should assume little about the database they occupy. Therefore, it's appropriate to write functions provided by an extension in a secure style that cannot be compromised by search-path-based attacks.

An extension that has the `superuser` property set to true must also consider security hazards for the actions taken within its installation and update scripts. It is not terribly difficult for a malicious user to create trojan-horse objects that will compromise later execution of a carelessly-written extension script, allowing that user to acquire superuser privileges.

If an extension is marked `trusted`, then its installation schema can be selected by the installing user, who might intentionally use an insecure schema in hopes of gaining superuser privileges. Therefore, a trusted extension is extremely exposed from a security standpoint, and all its script commands must be carefully examined to ensure that no compromise is possible.

Advice about writing functions securely is provided in [Section 39.17.6.1](#) below, and advice about writing installation scripts securely is provided in [Section 39.17.6.2](#).

39.17.6.1. Security Considerations for Extension Functions

SQL-language and PL-language functions provided by extensions are at risk of search-path-based attacks when they are executed, since parsing of these functions occurs at execution time not creation time.

The `CREATE FUNCTION` reference page contains advice about writing `SECURITY DEFINER` functions safely. It's good practice to apply those techniques for any function provided by an extension, since the function might be called by a high-privilege user.

If you cannot set the `search_path` to contain only secure schemas, assume that each unqualified name could resolve to an object that a malicious user has defined. Beware of constructs that depend on `search_path` implicitly; for example, `IN` and `CASE expression WHEN` always select an operator using the `search_path`. In their place, use `OPERATOR(schema.=) ANY` and `CASE WHEN expression`.

A general-purpose extension usually should not assume that it's been installed into a secure schema, which means that even schema-qualified references to its own objects are not entirely risk-free. For example, if the extension has defined a function `myschema.myfunc(bigint)` then a call such as `myschema.myfunc(42)` could be captured by a hostile function `myschema.myfunc(integer)`. Be careful that the data types of function and operator parameters exactly match the declared argument types, using explicit casts where necessary.

39.17.6.2. Security Considerations for Extension Scripts

An extension installation or update script should be written to guard against search-path-based attacks occurring when the script executes. If an object reference in the script can be made to resolve to some other object than the script author intended, then a compromise might occur immediately, or later when the mis-defined extension object is used.

DDL commands such as `CREATE FUNCTION` and `CREATE OPERATOR CLASS` are generally secure, but beware of any command having a general-purpose expression as a component. For example, `CREATE VIEW` needs to be vetted, as does a `DEFAULT` expression in `CREATE FUNCTION`.

Sometimes an extension script might need to execute general-purpose SQL, for example to make catalog adjustments that aren't possible via DDL. Be careful to execute such commands with a secure `search_path`; do *not* trust the path provided by `CREATE/ALTER EXTENSION` to be secure. Best practice is to temporarily set `search_path` to `pg_catalog`, `pg_temp` and insert references to the extension's installation schema explicitly where needed. (This practice might also be helpful for creating views.)

Cross-extension references are extremely difficult to make fully secure, partially because of uncertainty about which schema the other extension is in. The hazards are reduced if both extensions are installed in the same schema, because then a hostile object cannot be placed ahead of the referenced extension in the installation-time `search_path`. However, no mechanism currently exists to require that. For now, best practice is to not mark an extension trusted if it depends on another one, unless that other one is always installed in `pg_catalog`.

39.17.7. Extension Example

Here is a complete example of an SQL-only extension, a two-element composite type that can store any type of value in its slots, which are named “k” and “v”. Non-text values are automatically coerced to text for storage.

The script file `pair--1.0.sql` looks like this:

```
-- complain if script is sourced in psql, rather than via CREATE EXTENSION
\echo Use "CREATE EXTENSION pair" to load this file. \quit

CREATE TYPE pair AS ( k text, v text );

CREATE FUNCTION pair(text, text)
RETURNS pair LANGUAGE SQL AS 'SELECT ROW($1, $2)::@extschema@.pair;';
```

```
CREATE OPERATOR ~> (LEFTARG = text, RIGHTARG = text, FUNCTION = pair);

-- "SET search_path" is easy to get right, but qualified names perform better.
CREATE FUNCTION lower(pair)
RETURNS pair LANGUAGE SQL
AS 'SELECT ROW(lower($1.k), lower($1.v))::@extschema@.pair; '
SET search_path = pg_temp;

CREATE FUNCTION pair_concat(pair, pair)
RETURNS pair LANGUAGE SQL
AS 'SELECT ROW($1.k OPERATOR(pg_catalog.||) $2.k,
               $1.v OPERATOR(pg_catalog.||) $2.v)::@extschema@.pair;';
```

The control file `pair.control` looks like this:

```
# pair extension
comment = 'A key/value pair data type'
default_version = '1.0'
# cannot be relocatable because of use of @extschema@
relocatable = false
```

While you hardly need a makefile to install these two files into the correct directory, you could use a Makefile containing this:

```
EXTENSION = pair
DATA = pair--1.0.sql

PG_CONFIG = pg_config
PGXS := $(shell $(PG_CONFIG) --pgxs)
include $(PGXS)
```

This makefile relies on PGXS, which is described in [Section 39.18](#). The command `make install` will install the control and script files into the correct directory as reported by `pg_config`.

Once the files are installed, use the `CREATE EXTENSION` command to load the objects into any particular database.

39.18. Extension Building Infrastructure

If you are thinking about distributing your Postgres Pro extension modules, setting up a portable build system for them can be fairly difficult. Therefore the Postgres Pro installation provides a build infrastructure for extensions, called PGXS, so that simple extension modules can be built simply against an already installed server. PGXS is mainly intended for extensions that include C code, although it can be used for pure-SQL extensions too. Note that PGXS is not intended to be a universal build system framework that can be used to build any software interfacing to Postgres Pro; it simply automates common build rules for simple server extension modules. For more complicated packages, you might need to write your own build system.

To use the PGXS infrastructure for your extension, you must write a simple makefile. In the makefile, you need to set some variables and include the global PGXS makefile. Here is an example that builds an extension module named `isbn_issn`, consisting of a shared library containing some C code, an extension control file, an SQL script, an include file (only needed if other modules might need to access the extension functions without going via SQL), and a documentation text file:

```
MODULES = isbn_issn
EXTENSION = isbn_issn
DATA = isbn_issn--1.0.sql
DOCS = README.isbn_issn
HEADERS_isbn_issn = isbn_issn.h
```

```
PG_CONFIG = pg_config
PGXS := $(shell $(PG_CONFIG) --pgxs)
include $(PGXS)
```

The last three lines should always be the same. Earlier in the file, you assign variables or add custom make rules.

Set one of these three variables to specify what is built:

MODULES

list of shared-library objects to be built from source files with same stem (do not include library suffixes in this list)

MODULE_big

a shared library to build from multiple source files (list object files in OBJS)

PROGRAM

an executable program to build (list object files in OBJS)

The following variables can also be set:

EXTENSION

extension name(s); for each name you must provide an *extension.control* file, which will be installed into *prefix/share/extension*

MODULEDIR

subdirectory of *prefix/share* into which DATA and DOCS files should be installed (if not set, default is *extension* if EXTENSION is set, or *contrib* if not)

DATA

random files to install into *prefix/share/\$MODULEDIR*

DATA_built

random files to install into *prefix/share/\$MODULEDIR*, which need to be built first

DATA_TSEARCH

random files to install under *prefix/share/tsearch_data*

DOCS

random files to install under *prefix/doc/\$MODULEDIR*

HEADERS

HEADERS_built

Files to (optionally build and) install under *prefix/include/server/\$MODULEDIR/\$MODULE_big*.

Unlike DATA_built, files in HEADERS_built are not removed by the *clean* target; if you want them removed, also add them to EXTRA_CLEAN or add your own rules to do it.

HEADERS_\$MODULE

HEADERS_built_\$MODULE

Files to install (after building if specified) under *prefix/include/server/\$MODULEDIR/\$MODULE*, where \$MODULE must be a module name used in MODULES or MODULE_big.

Unlike DATA_built, files in HEADERS_built_\$MODULE are not removed by the *clean* target; if you want them removed, also add them to EXTRA_CLEAN or add your own rules to do it.

It is legal to use both variables for the same module, or any combination, unless you have two module names in the `MODULES` list that differ only by the presence of a prefix `built_`, which would cause ambiguity. In that (hopefully unlikely) case, you should use only the `HEADERS_built_$MODULE` variables.

`SCRIPTS`

script files (not binaries) to install into `prefix/bin`

`SCRIPTS_built`

script files (not binaries) to install into `prefix/bin`, which need to be built first

`REGRESS`

list of regression test cases (without suffix), see below

`REGRESS_OPTS`

additional switches to pass to `pg_regress`

`ISOLATION`

list of isolation test cases, see below for more details

`ISOLATION_OPTS`

additional switches to pass to `pg_isolation_regress`

`TAP_TESTS`

switch defining if TAP tests need to be run, see below

`NO_INSTALL`

don't define an `install` target, useful for test modules that don't need their build products to be installed

`NO_INSTALLCHECK`

don't define an `installcheck` target, useful e.g., if tests require special configuration, or don't use `pg_regress`

`EXTRA_CLEAN`

extra files to remove in `make clean`

`PG_CPPFLAGS`

will be prepended to `CPPFLAGS`

`PG_CFLAGS`

will be appended to `CFLAGS`

`PG_CXXFLAGS`

will be appended to `CXXFLAGS`

`PG_LDFLAGS`

will be prepended to `LDLIBS`

`PG_LIBS`

will be added to `PROGRAM` link line

`SHLIB_LINK`

will be added to `MODULE_big` link line

PG_CONFIG

path to `pg_config` program for the Postgres Pro installation to build against (typically just `pg_config` to use the first one in your `PATH`)

Put this makefile as `Makefile` in the directory which holds your extension. Then you can do `make` to compile, and then `make install` to install your module. By default, the extension is compiled and installed for the Postgres Pro installation that corresponds to the first `pg_config` program found in your `PATH`. You can use a different installation by setting `PG_CONFIG` to point to its `pg_config` program, either within the makefile or on the `make` command line.

You can also run `make` in a directory outside the source tree of your extension, if you want to keep the build directory separate. This procedure is also called a *VPATH* build. Here's how:

```
mkdir build_dir
cd build_dir
make -f /path/to/extension/source/tree/Makefile
make -f /path/to/extension/source/tree/Makefile install
```

Alternatively, you can set up a directory for a *VPATH* build in a similar way to how it is done for the core code. One way to do this is using the core script `config/prep_buildtree`. Once this has been done you can build by setting the `make` variable `VPATH` like this:

```
make VPATH=/path/to/extension/source/tree
make VPATH=/path/to/extension/source/tree install
```

This procedure can work with a greater variety of directory layouts.

The scripts listed in the `REGRESS` variable are used for regression testing of your module, which can be invoked by `make installcheck` after doing `make install`. For this to work you must have a running Postgres Pro server. The script files listed in `REGRESS` must appear in a subdirectory named `sql/` in your extension's directory. These files must have extension `.sql`, which must not be included in the `REGRESS` list in the makefile. For each test there should also be a file containing the expected output in a subdirectory named `expected/`, with the same stem and extension `.out`. `make installcheck` executes each test script with `psql`, and compares the resulting output to the matching expected file. Any differences will be written to the file `regression.diffs` in `diff -c` format. Note that trying to run a test that is missing its expected file will be reported as “trouble”, so make sure you have all expected files.

The scripts listed in the `ISOLATION` variable are used for tests stressing behavior of concurrent session with your module, which can be invoked by `make installcheck` after doing `make install`. For this to work you must have a running Postgres Pro server. The script files listed in `ISOLATION` must appear in a subdirectory named `specs/` in your extension's directory. These files must have extension `.spec`, which must not be included in the `ISOLATION` list in the makefile. For each test there should also be a file containing the expected output in a subdirectory named `expected/`, with the same stem and extension `.out`. `make installcheck` executes each test script, and compares the resulting output to the matching expected file. Any differences will be written to the file `output_iso/regression.diffs` in `diff -c` format. Note that trying to run a test that is missing its expected file will be reported as “trouble”, so make sure you have all expected files.

`TAP_TESTS` enables the use of TAP tests. Data from each run is present in a subdirectory named `tmp-check/`.

Tip

The easiest way to create the expected files is to create empty files, then do a test run (which will of course report differences). Inspect the actual result files found in the `results/` directory (for tests in `REGRESS`), or `output_iso/results/` directory (for tests in `ISOLATION`), then copy them to `expected/` if they match what you expect from the test.

Chapter 40. Triggers

This chapter provides general information about writing trigger functions. Trigger functions can be written in most of the available procedural languages, including PL/pgSQL ([Chapter 44](#)), PL/Tcl ([Chapter 45](#)), PL/Perl ([Chapter 46](#)), and PL/Python ([Chapter 47](#)). After reading this chapter, you should consult the chapter for your favorite procedural language to find out the language-specific details of writing a trigger in it.

It is also possible to write a trigger function in C, although most people find it easier to use one of the procedural languages. It is not currently possible to write a trigger function in the plain SQL function language.

40.1. Overview of Trigger Behavior

A trigger is a specification that the database should automatically execute a particular function whenever a certain type of operation is performed. Triggers can be attached to tables (partitioned or not), views, and foreign tables.

On tables and foreign tables, triggers can be defined to execute either before or after any `INSERT`, `UPDATE`, or `DELETE` operation, either once per modified row, or once per SQL statement. `UPDATE` triggers can moreover be set to fire only if certain columns are mentioned in the `SET` clause of the `UPDATE` statement. Triggers can also fire for `TRUNCATE` statements. If a trigger event occurs, the trigger's function is called at the appropriate time to handle the event.

On views, triggers can be defined to execute instead of `INSERT`, `UPDATE`, or `DELETE` operations. Such `INSTEAD OF` triggers are fired once for each row that needs to be modified in the view. It is the responsibility of the trigger's function to perform the necessary modifications to the view's underlying base table(s) and, where appropriate, return the modified row as it will appear in the view. Triggers on views can also be defined to execute once per SQL statement, before or after `INSERT`, `UPDATE`, or `DELETE` operations. However, such triggers are fired only if there is also an `INSTEAD OF` trigger on the view. Otherwise, any statement targeting the view must be rewritten into a statement affecting its underlying base table(s), and then the triggers that will be fired are the ones attached to the base table(s).

The trigger function must be defined before the trigger itself can be created. The trigger function must be declared as a function taking no arguments and returning type `trigger`. (The trigger function receives its input through a specially-passed `TriggerData` structure, not in the form of ordinary function arguments.)

Once a suitable trigger function has been created, the trigger is established with [CREATE TRIGGER](#). The same trigger function can be used for multiple triggers.

Postgres Pro offers both *per-row* triggers and *per-statement* triggers. With a per-row trigger, the trigger function is invoked once for each row that is affected by the statement that fired the trigger. In contrast, a per-statement trigger is invoked only once when an appropriate statement is executed, regardless of the number of rows affected by that statement. In particular, a statement that affects zero rows will still result in the execution of any applicable per-statement triggers. These two types of triggers are sometimes called *row-level* triggers and *statement-level* triggers, respectively. Triggers on `TRUNCATE` may only be defined at statement level, not per-row.

Triggers are also classified according to whether they fire *before*, *after*, or *instead of* the operation. These are referred to as `BEFORE` triggers, `AFTER` triggers, and `INSTEAD OF` triggers respectively. Statement-level `BEFORE` triggers naturally fire before the statement starts to do anything, while statement-level `AFTER` triggers fire at the very end of the statement. These types of triggers may be defined on tables, views, or foreign tables. Row-level `BEFORE` triggers fire immediately before a particular row is operated on, while row-level `AFTER` triggers fire at the end of the statement (but before any statement-level `AFTER` triggers). These types of triggers may only be defined on tables and foreign tables, not views. `INSTEAD OF` triggers may only be defined on views, and only at row level; they fire immediately as each row in the view is identified as needing to be operated on.

The execution of an `AFTER` trigger can be deferred to the end of the transaction, rather than the end of the statement, if it was defined as a *constraint trigger*. In all cases, a trigger is executed as part of the same transaction as the statement that triggered it, so if either the statement or the trigger causes an error, the effects of both will be rolled back.

If an `INSERT` contains an `ON CONFLICT DO UPDATE` clause, it is possible for row-level `BEFORE INSERT` and then `BEFORE UPDATE` triggers to be executed on triggered rows. Such interactions can be complex if the triggers are not idempotent because change made by `BEFORE INSERT` triggers will be seen by `BEFORE UPDATE` triggers, including changes to `EXCLUDED` columns.

Note that statement-level `UPDATE` triggers are executed when `ON CONFLICT DO UPDATE` is specified, regardless of whether or not any rows were affected by the `UPDATE` (and regardless of whether the alternative `UPDATE` path was ever taken). An `INSERT` with an `ON CONFLICT DO UPDATE` clause will execute statement-level `BEFORE INSERT` triggers first, then statement-level `BEFORE UPDATE` triggers, followed by statement-level `AFTER UPDATE` triggers and finally statement-level `AFTER INSERT` triggers.

A statement that targets a parent table in an inheritance or partitioning hierarchy does not cause the statement-level triggers of affected child tables to be fired; only the parent table's statement-level triggers are fired. However, row-level triggers of any affected child tables will be fired.

If an `UPDATE` on a partitioned table causes a row to move to another partition, it will be performed as a `DELETE` from the original partition followed by an `INSERT` into the new partition. In this case, all row-level `BEFORE UPDATE` triggers and all row-level `BEFORE DELETE` triggers are fired on the original partition. Then all row-level `BEFORE INSERT` triggers are fired on the destination partition. The possibility of surprising outcomes should be considered when all these triggers affect the row being moved. As far as `AFTER ROW` triggers are concerned, `AFTER DELETE` and `AFTER INSERT` triggers are applied; but `AFTER UPDATE` triggers are not applied because the `UPDATE` has been converted to a `DELETE` and an `INSERT`. As far as statement-level triggers are concerned, none of the `DELETE` or `INSERT` triggers are fired, even if row movement occurs; only the `UPDATE` triggers defined on the target table used in the `UPDATE` statement will be fired.

No separate triggers are defined for `MERGE`. Instead, statement-level or row-level `UPDATE`, `DELETE`, and `INSERT` triggers are fired depending on (for statement-level triggers) what actions are specified in the `MERGE` query and (for row-level triggers) what actions are performed.

While running a `MERGE` command, statement-level `BEFORE` and `AFTER` triggers are fired for events specified in the actions of the `MERGE` command, irrespective of whether or not the action is ultimately performed. This is the same as an `UPDATE` statement that updates no rows, yet statement-level triggers are fired. The row-level triggers are fired only when a row is actually updated, inserted or deleted. So it's perfectly legal that while statement-level triggers are fired for certain types of action, no row-level triggers are fired for the same kind of action.

Trigger functions invoked by per-statement triggers should always return `NULL`. Trigger functions invoked by per-row triggers can return a table row (a value of type `HeapTuple`) to the calling executor, if they choose. A row-level trigger fired before an operation has the following choices:

- It can return `NULL` to skip the operation for the current row. This instructs the executor to not perform the row-level operation that invoked the trigger (the insertion, modification, or deletion of a particular table row).
- For row-level `INSERT` and `UPDATE` triggers only, the returned row becomes the row that will be inserted or will replace the row being updated. This allows the trigger function to modify the row being inserted or updated.

A row-level `BEFORE` trigger that does not intend to cause either of these behaviors must be careful to return as its result the same row that was passed in (that is, the `NEW` row for `INSERT` and `UPDATE` triggers, the `OLD` row for `DELETE` triggers).

A row-level `INSTEAD OF` trigger should either return `NULL` to indicate that it did not modify any data from the view's underlying base tables, or it should return the view row that was passed in (the `NEW` row for `INSERT` and `UPDATE` operations, or the `OLD` row for `DELETE` operations). A nonnull return value is used to

signal that the trigger performed the necessary data modifications in the view. This will cause the count of the number of rows affected by the command to be incremented. For `INSERT` and `UPDATE` operations only, the trigger may modify the `NEW` row before returning it. This will change the data returned by `INSERT RETURNING` or `UPDATE RETURNING`, and is useful when the view will not show exactly the same data that was provided.

The return value is ignored for row-level triggers fired after an operation, and so they can return `NULL`.

Some considerations apply for generated columns. Stored generated columns are computed after `BEFORE` triggers and before `AFTER` triggers. Therefore, the generated value can be inspected in `AFTER` triggers. In `BEFORE` triggers, the `OLD` row contains the old generated value, as one would expect, but the `NEW` row does not yet contain the new generated value and should not be accessed. In the C language interface, the content of the column is undefined at this point; a higher-level programming language should prevent access to a stored generated column in the `NEW` row in a `BEFORE` trigger. Changes to the value of a generated column in a `BEFORE` trigger are ignored and will be overwritten.

If more than one trigger is defined for the same event on the same relation, the triggers will be fired in alphabetical order by trigger name. In the case of `BEFORE` and `INSTEAD OF` triggers, the possibly-modified row returned by each trigger becomes the input to the next trigger. If any `BEFORE` or `INSTEAD OF` trigger returns `NULL`, the operation is abandoned for that row and subsequent triggers are not fired (for that row).

A trigger definition can also specify a Boolean `WHEN` condition, which will be tested to see whether the trigger should be fired. In row-level triggers the `WHEN` condition can examine the old and/or new values of columns of the row. (Statement-level triggers can also have `WHEN` conditions, although the feature is not so useful for them.) In a `BEFORE` trigger, the `WHEN` condition is evaluated just before the function is or would be executed, so using `WHEN` is not materially different from testing the same condition at the beginning of the trigger function. However, in an `AFTER` trigger, the `WHEN` condition is evaluated just after the row update occurs, and it determines whether an event is queued to fire the trigger at the end of statement. So when an `AFTER` trigger's `WHEN` condition does not return true, it is not necessary to queue an event nor to re-fetch the row at end of statement. This can result in significant speedups in statements that modify many rows, if the trigger only needs to be fired for a few of the rows. `INSTEAD OF` triggers do not support `WHEN` conditions.

Typically, row-level `BEFORE` triggers are used for checking or modifying the data that will be inserted or updated. For example, a `BEFORE` trigger might be used to insert the current time into a `timestamp` column, or to check that two elements of the row are consistent. Row-level `AFTER` triggers are most sensibly used to propagate the updates to other tables, or make consistency checks against other tables. The reason for this division of labor is that an `AFTER` trigger can be certain it is seeing the final value of the row, while a `BEFORE` trigger cannot; there might be other `BEFORE` triggers firing after it. If you have no specific reason to make a trigger `BEFORE` or `AFTER`, the `BEFORE` case is more efficient, since the information about the operation doesn't have to be saved until end of statement.

If a trigger function executes SQL commands then these commands might fire triggers again. This is known as cascading triggers. There is no direct limitation on the number of cascade levels. It is possible for cascades to cause a recursive invocation of the same trigger; for example, an `INSERT` trigger might execute a command that inserts an additional row into the same table, causing the `INSERT` trigger to be fired again. It is the trigger programmer's responsibility to avoid infinite recursion in such scenarios.

If a foreign key constraint specifies referential actions (that is, cascading updates or deletes), those actions are performed via ordinary SQL update or delete commands on the referencing table. In particular, any triggers that exist on the referencing table will be fired for those changes. If such a trigger modifies or blocks the effect of one of these commands, the end result could be to break referential integrity. It is the trigger programmer's responsibility to avoid that.

When a trigger is being defined, arguments can be specified for it. The purpose of including arguments in the trigger definition is to allow different triggers with similar requirements to call the same function. As an example, there could be a generalized trigger function that takes as its arguments two column names and puts the current user in one and the current time stamp in the other. Properly written, this

trigger function would be independent of the specific table it is triggering on. So the same function could be used for `INSERT` events on any table with suitable columns, to automatically track creation of records in a transaction table for example. It could also be used to track last-update events if defined as an `UPDATE` trigger.

Each programming language that supports triggers has its own method for making the trigger input data available to the trigger function. This input data includes the type of trigger event (e.g., `INSERT` or `UPDATE`) as well as any arguments that were listed in `CREATE TRIGGER`. For a row-level trigger, the input data also includes the `NEW` row for `INSERT` and `UPDATE` triggers, and/or the `OLD` row for `UPDATE` and `DELETE` triggers.

By default, statement-level triggers do not have any way to examine the individual row(s) modified by the statement. But an `AFTER STATEMENT` trigger can request that *transition tables* be created to make the sets of affected rows available to the trigger. `AFTER ROW` triggers can also request transition tables, so that they can see the total changes in the table as well as the change in the individual row they are currently being fired for. The method for examining the transition tables again depends on the programming language that is being used, but the typical approach is to make the transition tables act like read-only temporary tables that can be accessed by SQL commands issued within the trigger function.

40.2. Visibility of Data Changes

If you execute SQL commands in your trigger function, and these commands access the table that the trigger is for, then you need to be aware of the data visibility rules, because they determine whether these SQL commands will see the data change that the trigger is fired for. Briefly:

- Statement-level triggers follow simple visibility rules: none of the changes made by a statement are visible to statement-level `BEFORE` triggers, whereas all modifications are visible to statement-level `AFTER` triggers.
- The data change (insertion, update, or deletion) causing the trigger to fire is naturally *not* visible to SQL commands executed in a row-level `BEFORE` trigger, because it hasn't happened yet.
- However, SQL commands executed in a row-level `BEFORE` trigger *will* see the effects of data changes for rows previously processed in the same outer command. This requires caution, since the ordering of these change events is not in general predictable; an SQL command that affects multiple rows can visit the rows in any order.
- Similarly, a row-level `INSTEAD OF` trigger will see the effects of data changes made by previous firings of `INSTEAD OF` triggers in the same outer command.
- When a row-level `AFTER` trigger is fired, all data changes made by the outer command are already complete, and are visible to the invoked trigger function.

If your trigger function is written in any of the standard procedural languages, then the above statements apply only if the function is declared `VOLATILE`. Functions that are declared `STABLE` or `IMMUTABLE` will not see changes made by the calling command in any case.

Further information about data visibility rules can be found in [Section 48.5](#). The example in [Section 40.4](#) contains a demonstration of these rules.

40.3. Writing Trigger Functions in C

This section describes the low-level details of the interface to a trigger function. This information is only needed when writing trigger functions in C. If you are using a higher-level language then these details are handled for you. In most cases you should consider using a procedural language before writing your triggers in C. The documentation of each procedural language explains how to write a trigger in that language.

Trigger functions must use the “version 1” function manager interface.

When a function is called by the trigger manager, it is not passed any normal arguments, but it is passed a “context” pointer pointing to a `TriggerData` structure. C functions can check whether they were called from the trigger manager or not by executing the macro:

```
CALLED_AS_TRIGGER(fcinfo)
```

which expands to:

```
((fcinfo)->context != NULL && IsA((fcinfo)->context, TriggerData))
```

If this returns true, then it is safe to cast `fcinfo->context` to type `TriggerData *` and make use of the pointed-to `TriggerData` structure. The function must *not* alter the `TriggerData` structure or any of the data it points to.

struct `TriggerData` is defined in `commands/trigger.h`:

```
typedef struct TriggerData
{
    NodeTag          type;
    TriggerEvent      tg_event;
    Relation          tg_relation;
    HeapTuple         tg_trigtuple;
    HeapTuple         tg_newtuple;
    Trigger           *tg_trigger;
    TupleTableSlot    *tg_trigslot;
    TupleTableSlot    *tg_newslot;
    Tuplestorestate   *tg_oldtable;
    Tuplestorestate   *tg_newtable;
    const Bitmapset   *tg_updatedcols;
} TriggerData;
```

where the members are defined as follows:

type

Always `T_TriggerData`.

tg_event

Describes the event for which the function is called. You can use the following macros to examine `tg_event`:

`TRIGGER_FIRED_BEFORE(tg_event)`

Returns true if the trigger fired before the operation.

`TRIGGER_FIRED_AFTER(tg_event)`

Returns true if the trigger fired after the operation.

`TRIGGER_FIRED_INSTEAD(tg_event)`

Returns true if the trigger fired instead of the operation.

`TRIGGER_FIRED_FOR_ROW(tg_event)`

Returns true if the trigger fired for a row-level event.

`TRIGGER_FIRED_FOR_STATEMENT(tg_event)`

Returns true if the trigger fired for a statement-level event.

`TRIGGER_FIRED_BY_INSERT(tg_event)`

Returns true if the trigger was fired by an `INSERT` command.

`TRIGGER_FIRED_BY_UPDATE(tg_event)`

Returns true if the trigger was fired by an `UPDATE` command.

`TRIGGER_FIRED_BY_DELETE(tg_event)`

Returns true if the trigger was fired by a `DELETE` command.

TRIGGER_FIRED_BY_TRUNCATE(tg_event)

Returns true if the trigger was fired by a TRUNCATE command.

tg_relation

A pointer to a structure describing the relation that the trigger fired for. Look at `utils/rel.h` for details about this structure. The most interesting things are `tg_relation->rd_att` (descriptor of the relation tuples) and `tg_relation->rd_rel->relname` (relation name; the type is not `char*` but `NameData`; use `SPI_getrelname(tg_relation)` to get a `char*` if you need a copy of the name).

tg_trigtuple

A pointer to the row for which the trigger was fired. This is the row being inserted, updated, or deleted. If this trigger was fired for an INSERT or DELETE then this is what you should return from the function if you don't want to replace the row with a different one (in the case of INSERT) or skip the operation. For triggers on foreign tables, values of system columns herein are unspecified.

tg_newtuple

A pointer to the new version of the row, if the trigger was fired for an UPDATE, and NULL if it is for an INSERT or a DELETE. This is what you have to return from the function if the event is an UPDATE and you don't want to replace this row by a different one or skip the operation. For triggers on foreign tables, values of system columns herein are unspecified.

tg_trigger

A pointer to a structure of type `Trigger`, defined in `utils/reltrigger.h`:

```
typedef struct Trigger
{
    Oid          tgoid;
    char         *tgname;
    Oid          tgfoid;
    int16        tgtype;
    char         tgenabled;
    bool         tgisinternal;
    bool         tgisclone;
    Oid          tgconstrrelid;
    Oid          tgconstrindid;
    Oid          tgconstraint;
    bool         tgdeferrable;
    bool         tginitdeferred;
    int16        tgnargs;
    int16        tgnattr;
    int16        *tgattr;
    char         **tgargs;
    char         *tgqual;
    char         *tgoldtable;
    char         *tgnewtable;
} Trigger;
```

where `tgname` is the trigger's name, `tgnargs` is the number of arguments in `tgargs`, and `tgargs` is an array of pointers to the arguments specified in the CREATE TRIGGER statement. The other members are for internal use only.

tg_trigslot

The slot containing `tg_trigtuple`, or a NULL pointer if there is no such tuple.

tg_newslot

The slot containing `tg_newtuple`, or a NULL pointer if there is no such tuple.

tg_oldtable

A pointer to a structure of type `Tuplestorestate` containing zero or more rows in the format specified by `tg_relation`, or a `NULL` pointer if there is no `OLD TABLE` transition relation.

tg_newtable

A pointer to a structure of type `Tuplestorestate` containing zero or more rows in the format specified by `tg_relation`, or a `NULL` pointer if there is no `NEW TABLE` transition relation.

tg_updatedcols

For `UPDATE` triggers, a bitmap set indicating the columns that were updated by the triggering command. Generic trigger functions can use this to optimize actions by not having to deal with columns that were not changed.

As an example, to determine whether a column with attribute number `attnum` (1-based) is a member of this bitmap set, call `bms_is_member(attnum - FirstLowInvalidHeapAttributeNumber, trigdata->tg_updatedcols)`.

For triggers other than `UPDATE` triggers, this will be `NULL`.

To allow queries issued through SPI to reference transition tables, see [SPI_register_trigger_data](#).

A trigger function must return either a `HeapTuple` pointer or a `NULL` pointer (*not* an SQL null value, that is, do not set `isNull` true). Be careful to return either `tg_trigtuple` or `tg_newtuple`, as appropriate, if you don't want to modify the row being operated on.

40.4. A Complete Trigger Example

Here is a very simple example of a trigger function written in C. (Examples of triggers written in procedural languages can be found in the documentation of the procedural languages.)

The function `trigf` reports the number of rows in the table `ttest` and skips the actual operation if the command attempts to insert a null value into the column `x`. (So the trigger acts as a not-null constraint but doesn't abort the transaction.)

First, the table definition:

```
CREATE TABLE ttest (  
    x integer  
);
```

This is the source code of the trigger function:

```
#include "postgres.h"  
#include "fmgr.h"  
#include "executor/spi.h"      /* this is what you need to work with SPI */  
#include "commands/trigger.h" /* ... triggers ... */  
#include "utils/rel.h"        /* ... and relations */  
  
PG_MODULE_MAGIC;  
  
PG_FUNCTION_INFO_V1(trigf);  
  
Datum  
trigf(PG_FUNCTION_ARGS)  
{  
    TriggerData *trigdata = (TriggerData *) fcinfo->context;  
    TupleDesc   tupdesc;  
    HeapTuple   rettuple;  
    char        *when;  
    bool        checknull = false;
```

```
bool        isnull;
int         ret, i;

/* make sure it's called as a trigger at all */
if (!CALLED_AS_TRIGGER(fcinfo))
    elog(ERROR, "trigf: not called by trigger manager");

/* tuple to return to executor */
if (TRIGGER_FIRED_BY_UPDATE(trigdata->tg_event))
    rettupple = trigdata->tg_newtuple;
else
    rettupple = trigdata->tg_trigtuple;

/* check for null values */
if (!TRIGGER_FIRED_BY_DELETE(trigdata->tg_event)
    && TRIGGER_FIRED_BEFORE(trigdata->tg_event))
    checknull = true;

if (TRIGGER_FIRED_BEFORE(trigdata->tg_event))
    when = "before";
else
    when = "after ";

tupdesc = trigdata->tg_relation->rd_att;

/* connect to SPI manager */
if ((ret = SPI_connect()) < 0)
    elog(ERROR, "trigf (fired %s): SPI_connect returned %d", when, ret);

/* get number of rows in table */
ret = SPI_exec("SELECT count(*) FROM ttest", 0);

if (ret < 0)
    elog(ERROR, "trigf (fired %s): SPI_exec returned %d", when, ret);

/* count(*) returns int8, so be careful to convert */
i = DatumGetInt64(SPI_getbinval(SPI_tuptable->vals[0],
                                SPI_tuptable->tupdesc,
                                1,
                                &isnull));

elog (INFO, "trigf (fired %s): there are %d rows in ttest", when, i);

SPI_finish();

if (checknull)
{
    SPI_getbinval(rettuple, tupdesc, 1, &isnull);
    if (isnull)
        rettupple = NULL;
}

return PointerGetDatum(rettuple);
}
```

After you have compiled the source code (see [Section 39.10.5](#)), declare the function and the triggers:

```
CREATE FUNCTION trigf() RETURNS trigger
```



```
AS 'filename'
LANGUAGE C;
```

```
CREATE TRIGGER tbefore BEFORE INSERT OR UPDATE OR DELETE ON ttest
FOR EACH ROW EXECUTE FUNCTION trigf();
```

```
CREATE TRIGGER tafter AFTER INSERT OR UPDATE OR DELETE ON ttest
FOR EACH ROW EXECUTE FUNCTION trigf();
```

Now you can test the operation of the trigger:

```
=> INSERT INTO ttest VALUES (NULL);
INFO:  trigf (fired before): there are 0 rows in ttest
INSERT 0 0
```

```
-- Insertion skipped and AFTER trigger is not fired
```

```
=> SELECT * FROM ttest;
 x
---
(0 rows)
```

```
=> INSERT INTO ttest VALUES (1);
INFO:  trigf (fired before): there are 0 rows in ttest
INFO:  trigf (fired after ): there are 1 rows in ttest
      ^^^^^^^^
```

remember what we said about visibility.

```
INSERT 167793 1
vac=> SELECT * FROM ttest;
 x
---
 1
(1 row)
```

```
=> INSERT INTO ttest SELECT x * 2 FROM ttest;
INFO:  trigf (fired before): there are 1 rows in ttest
INFO:  trigf (fired after ): there are 2 rows in ttest
      ^^^^^^
```

remember what we said about visibility.

```
INSERT 167794 1
=> SELECT * FROM ttest;
 x
---
 1
 2
(2 rows)
```

```
=> UPDATE ttest SET x = NULL WHERE x = 2;
INFO:  trigf (fired before): there are 2 rows in ttest
UPDATE 0
```

```
=> UPDATE ttest SET x = 4 WHERE x = 2;
INFO:  trigf (fired before): there are 2 rows in ttest
INFO:  trigf (fired after ): there are 2 rows in ttest
UPDATE 1
```

```
vac=> SELECT * FROM ttest;
 x
---
 1
 4
```

(2 rows)

```
=> DELETE FROM ttest;
```

```
INFO:  trigf (fired before): there are 2 rows in ttest
```

```
INFO:  trigf (fired before): there are 1 rows in ttest
```

```
INFO:  trigf (fired after ) : there are 0 rows in ttest
```

```
INFO:  trigf (fired after ) : there are 0 rows in ttest
```

^^^^^

remember what we said about visibility.

```
DELETE 2
```

```
=> SELECT * FROM ttest;
```

```
 x
```

```
---
```

(0 rows)

Chapter 41. Event Triggers

To supplement the trigger mechanism discussed in [Chapter 40](#), Postgres Pro also provides event triggers. Unlike regular triggers, which are attached to a single table and capture only DML events, event triggers are global to a particular database and are capable of capturing DDL events.

Like regular triggers, event triggers can be written in any procedural language that includes event trigger support, or in C, but not in plain SQL.

41.1. Overview of Event Trigger Behavior

An event trigger fires whenever the event with which it is associated occurs in the database in which it is defined. Currently, the only supported events are `login`, `ddl_command_start`, `ddl_command_end`, `table_rewrite` and `sql_drop`. Support for additional events may be added in future releases.

The `login` event occurs when an authenticated user logs into the system. Any bug in a trigger procedure for this event may prevent successful login to the system. Such bugs may be worked around by setting [event_triggers](#) to `false` either in a connection string or configuration file. Alternatively, you can restart the system in single-user mode (as event triggers are disabled in this mode). See the [postgres](#) reference page for details about using single-user mode. The `login` event will also fire on standby servers. To prevent servers from becoming inaccessible, such triggers must avoid writing anything to the database when running on a standby. Also, it's recommended to avoid long-running queries in `login` event triggers. Note that, for instance, canceling a connection in `psql` will not cancel the in-progress `login` trigger.

The `ddl_command_start` event occurs just before the execution of a `CREATE`, `ALTER`, `DROP`, `SECURITY LABEL`, `COMMENT`, `GRANT` or `REVOKE` command. No check whether the affected object exists or doesn't exist is performed before the event trigger fires. As an exception, however, this event does not occur for DDL commands targeting shared objects — databases, roles, and tablespaces — or for commands targeting event triggers themselves. The event trigger mechanism does not support these object types. `ddl_command_start` also occurs just before the execution of a `SELECT INTO` command, since this is equivalent to `CREATE TABLE AS`.

The `ddl_command_end` event occurs just after the execution of this same set of commands. To obtain more details on the DDL operations that took place, use the set-returning function `pg_event_trigger_ddl_commands()` from the `ddl_command_end` event trigger code (see [Section 9.30](#)). Note that the trigger fires after the actions have taken place (but before the transaction commits), and thus the system catalogs can be read as already changed.

The `sql_drop` event occurs just before the `ddl_command_end` event trigger for any operation that drops database objects. To list the objects that have been dropped, use the set-returning function `pg_event_trigger_dropped_objects()` from the `sql_drop` event trigger code (see [Section 9.30](#)). Note that the trigger is executed after the objects have been deleted from the system catalogs, so it's not possible to look them up anymore.

The `table_rewrite` event occurs just before a table is rewritten by some actions of the commands `ALTER TABLE` and `ALTER TYPE`. While other control statements are available to rewrite a table, like `CLUSTER` and `VACUUM`, the `table_rewrite` event is not triggered by them. To find the OID of the table that was rewritten, use the function `pg_event_trigger_table_rewrite_oid()` (see [Section 9.30](#)). To discover the reason(s) for the rewrite, use the function `pg_event_trigger_table_rewrite_reason()`.

Event triggers (like other functions) cannot be executed in an aborted transaction. Thus, if a DDL command fails with an error, any associated `ddl_command_end` triggers will not be executed. Conversely, if a `ddl_command_start` trigger fails with an error, no further event triggers will fire, and no attempt will be made to execute the command itself. Similarly, if a `ddl_command_end` trigger fails with an error, the effects of the DDL statement will be rolled back, just as they would be in any other case where the containing transaction aborts.

For a complete list of commands supported by the event trigger mechanism, see [Section 41.2](#).

Event triggers are created using the command **CREATE EVENT TRIGGER**. In order to create an event trigger, you must first create a function with the special return type `event_trigger`. This function need not (and may not) return a value; the return type serves merely as a signal that the function is to be invoked as an event trigger.

If more than one event trigger is defined for a particular event, they will fire in alphabetical order by trigger name.

A trigger definition can also specify a **WHEN** condition so that, for example, a `ddl_command_start` trigger can be fired only for particular commands which the user wishes to intercept. A common use of such triggers is to restrict the range of DDL operations which users may perform.

41.2. Event Trigger Firing Matrix

Table 41.1 lists all commands for which event triggers are supported.

Table 41.1. Event Trigger Support by Command Tag

Command Tag	ddl_com- mand_ start	ddl_com- mand_end	sql_drop	table_ rewrite	Notes
ALTER AGGREGATE	X	X	–	–	
ALTER COLLATION	X	X	–	–	
ALTER CONVERSION	X	X	–	–	
ALTER DOMAIN	X	X	–	–	
ALTER DEFAULT PRIVILEGES	X	X	–	–	
ALTER EXTENSION	X	X	–	–	
ALTER FOREIGN DATA WRAP- PER	X	X	–	–	
ALTER FOREIGN TABLE	X	X	X	–	
ALTER FUNCTION	X	X	–	–	
ALTER LANGUAGE	X	X	–	–	
ALTER LARGE OBJECT	X	X	–	–	
ALTER MATERIALIZED VIEW	X	X	–	X	
ALTER OPERATOR	X	X	–	–	
ALTER OPERATOR CLASS	X	X	–	–	
ALTER OPERATOR FAMILY	X	X	–	–	
ALTER POLICY	X	X	–	–	
ALTER PROCEDURE	X	X	–	–	
ALTER PUBLICATION	X	X	–	–	
ALTER ROUTINE	X	X	–	–	
ALTER SCHEMA	X	X	–	–	
ALTER SEQUENCE	X	X	–	–	
ALTER SERVER	X	X	–	–	
ALTER STATISTICS	X	X	–	–	
ALTER SUBSCRIPTION	X	X	–	–	
ALTER TABLE	X	X	X	X	

Event Triggers

Command Tag	ddl_com- mand_ start	ddl_com- mand_end	sql_drop	table_ rewrite	Notes
ALTER TEXT SEARCH CON- FIGURATION	X	X	–	–	
ALTER TEXT SEARCH DIC- TIONARY	X	X	–	–	
ALTER TEXT SEARCH PARSER	X	X	–	–	
ALTER TEXT SEARCH TEM- PLATE	X	X	–	–	
ALTER TRIGGER	X	X	–	–	
ALTER TYPE	X	X	–	X	
ALTER USER MAPPING	X	X	–	–	
ALTER VIEW	X	X	–	–	
COMMENT	X	X	–	–	Only for local objects
CREATE ACCESS METHOD	X	X	–	–	
CREATE AGGREGATE	X	X	–	–	
CREATE CAST	X	X	–	–	
CREATE COLLATION	X	X	–	–	
CREATE CONVERSION	X	X	–	–	
CREATE DOMAIN	X	X	–	–	
CREATE EXTENSION	X	X	–	–	
CREATE FOREIGN DATA WRAPPER	X	X	–	–	
CREATE FOREIGN TABLE	X	X	–	–	
CREATE FUNCTION	X	X	–	–	
CREATE INDEX	X	X	–	–	
CREATE LANGUAGE	X	X	–	–	
CREATE MATERIALIZED VIEW	X	X	–	–	
CREATE OPERATOR	X	X	–	–	
CREATE OPERATOR CLASS	X	X	–	–	
CREATE OPERATOR FAMILY	X	X	–	–	
CREATE POLICY	X	X	–	–	
CREATE PROCEDURE	X	X	–	–	
CREATE PUBLICATION	X	X	–	–	
CREATE RULE	X	X	–	–	
CREATE SCHEMA	X	X	–	–	
CREATE SEQUENCE	X	X	–	–	
CREATE SERVER	X	X	–	–	
CREATE STATISTICS	X	X	–	–	
CREATE SUBSCRIPTION	X	X	–	–	
CREATE TABLE	X	X	–	–	

Event Triggers

Command Tag	ddl_com- mand_ start	ddl_com- mand_end	sql_drop	table_ rewrite	Notes
CREATE TABLE AS	X	X	–	–	
CREATE TEXT SEARCH CON- FIGURATION	X	X	–	–	
CREATE TEXT SEARCH DIC- TIONARY	X	X	–	–	
CREATE TEXT SEARCH PARSER	X	X	–	–	
CREATE TEXT SEARCH TEM- PLATE	X	X	–	–	
CREATE TRIGGER	X	X	–	–	
CREATE TYPE	X	X	–	–	
CREATE USER MAPPING	X	X	–	–	
CREATE VIEW	X	X	–	–	
DROP ACCESS METHOD	X	X	X	–	
DROP AGGREGATE	X	X	X	–	
DROP CAST	X	X	X	–	
DROP COLLATION	X	X	X	–	
DROP CONVERSION	X	X	X	–	
DROP DOMAIN	X	X	X	–	
DROP EXTENSION	X	X	X	–	
DROP FOREIGN DATA WRAP- PER	X	X	X	–	
DROP FOREIGN TABLE	X	X	X	–	
DROP FUNCTION	X	X	X	–	
DROP INDEX	X	X	X	–	
DROP LANGUAGE	X	X	X	–	
DROP MATERIALIZED VIEW	X	X	X	–	
DROP OPERATOR	X	X	X	–	
DROP OPERATOR CLASS	X	X	X	–	
DROP OPERATOR FAMILY	X	X	X	–	
DROP OWNED	X	X	X	–	
DROP POLICY	X	X	X	–	
DROP PROCEDURE	X	X	X	–	
DROP PUBLICATION	X	X	X	–	
DROP ROUTINE	X	X	X	–	
DROP RULE	X	X	X	–	
DROP SCHEMA	X	X	X	–	
DROP SEQUENCE	X	X	X	–	
DROP SERVER	X	X	X	–	
DROP STATISTICS	X	X	X	–	

Command Tag	ddl_com- mand_ start	ddl_com- mand_end	sql_drop	table_ rewrite	Notes
DROP SUBSCRIPTION	X	X	X	–	
DROP TABLE	X	X	X	–	
DROP TEXT SEARCH CON- FIGURATION	X	X	X	–	
DROP TEXT SEARCH DIC- TIONARY	X	X	X	–	
DROP TEXT SEARCH PARSER	X	X	X	–	
DROP TEXT SEARCH TEM- PLATE	X	X	X	–	
DROP TRIGGER	X	X	X	–	
DROP TYPE	X	X	X	–	
DROP USER MAPPING	X	X	X	–	
DROP VIEW	X	X	X	–	
GRANT	X	X	–	–	Only for local objects
IMPORT FOREIGN SCHEMA	X	X	–	–	
REFRESH MATERIALIZED VIEW	X	X	–	–	
REINDEX	X	X	–	–	
REVOKE	X	X	–	–	Only for local objects
SECURITY LABEL	X	X	–	–	Only for local objects
SELECT INTO	X	X	–	–	

41.3. Writing Event Trigger Functions in C

This section describes the low-level details of the interface to an event trigger function. This information is only needed when writing event trigger functions in C. If you are using a higher-level language then these details are handled for you. In most cases you should consider using a procedural language before writing your event triggers in C. The documentation of each procedural language explains how to write an event trigger in that language.

Event trigger functions must use the “version 1” function manager interface.

When a function is called by the event trigger manager, it is not passed any normal arguments, but it is passed a “context” pointer pointing to a `EventTriggerData` structure. C functions can check whether they were called from the event trigger manager or not by executing the macro:

```
CALLED_AS_EVENT_TRIGGER(fcinfo)
```

which expands to:

```
((fcinfo)->context != NULL && IsA((fcinfo)->context, EventTriggerData))
```

If this returns true, then it is safe to cast `fcinfo->context` to type `EventTriggerData *` and make use of the pointed-to `EventTriggerData` structure. The function must *not* alter the `EventTriggerData` structure or any of the data it points to.

`struct EventTriggerData` is defined in `commands/event_trigger.h`:

```
typedef struct EventTriggerData
{
    NodeTag      type;
    const char *event;      /* event name */
    Node         *parsetree; /* parse tree */
    CommandTag    tag;      /* command tag */
} EventTriggerData;
```

where the members are defined as follows:

type

Always `T_EventTriggerData`.

event

Describes the event for which the function is called, one of "login", "ddl_command_start", "ddl_command_end", "sql_drop", "table_rewrite". See [Section 41.1](#) for the meaning of these events.

parsetree

A pointer to the parse tree of the command. Check the Postgres Pro source code for details. The parse tree structure is subject to change without notice.

tag

The command tag associated with the event for which the event trigger is run, for example "CREATE FUNCTION".

An event trigger function must return a `NULL` pointer (*not* an SQL null value, that is, do not set `isNull` true).

41.4. A Complete Event Trigger Example

Here is a very simple example of an event trigger function written in C. (Examples of triggers written in procedural languages can be found in the documentation of the procedural languages.)

The function `nodd1` raises an exception each time it is called. The event trigger definition associated the function with the `ddl_command_start` event. The effect is that all DDL commands (with the exceptions mentioned in [Section 41.1](#)) are prevented from running.

This is the source code of the trigger function:

```
#include "postgres.h"

#include "commands/event_trigger.h"
#include "fmgr.h"

PG_MODULE_MAGIC;

PG_FUNCTION_INFO_V1(nodd1);

Datum
nodd1(PG_FUNCTION_ARGS)
{
    EventTriggerData *trigdata;

    if (!CALLED_AS_EVENT_TRIGGER(fcinfo)) /* internal error */
        elog(ERROR, "not fired by event trigger manager");

    trigdata = (EventTriggerData *) fcinfo->context;
```



```
ereport (ERROR,
        (errcode(ERRCODE_INSUFFICIENT_PRIVILEGE),
          errmsg("command \"%s\" denied",
                GetCommandTagName(trigdata->tag))));

PG_RETURN_NULL();
}
```

After you have compiled the source code (see [Section 39.10.5](#)), declare the function and the triggers:

```
CREATE FUNCTION nodd1() RETURNS event_trigger
AS 'nodd1' LANGUAGE C;

CREATE EVENT TRIGGER nodd1 ON ddl_command_start
EXECUTE FUNCTION nodd1();
```

Now you can test the operation of the trigger:

```
=# \dy
                                List of event triggers
 Name |          Event          | Owner | Enabled | Function | Tags
-----+-----+-----+-----+-----+-----
 nodd1 | ddl_command_start | dim   | enabled | nodd1    |
(1 row)

=# CREATE TABLE foo(id serial);
ERROR:  command "CREATE TABLE" denied
```

In this situation, in order to be able to run some DDL commands when you need to do so, you have to either drop the event trigger or disable it. It can be convenient to disable the trigger for only the duration of a transaction:

```
BEGIN;
ALTER EVENT TRIGGER nodd1 DISABLE;
CREATE TABLE foo (id serial);
ALTER EVENT TRIGGER nodd1 ENABLE;
COMMIT;
```

(Recall that DDL commands on event triggers themselves are not affected by event triggers.)

41.5. A Table Rewrite Event Trigger Example

Thanks to the `table_rewrite` event, it is possible to implement a table rewriting policy only allowing the rewrite in maintenance windows.

Here's an example implementing such a policy.

```
CREATE OR REPLACE FUNCTION no_rewrite()
RETURNS event_trigger
LANGUAGE plpgsql AS
$$
---
--- Implement local Table Rewriting policy:
---   public.foo is not allowed rewriting, ever
---   other tables are only allowed rewriting between 1am and 6am
---   unless they have more than 100 blocks
---
DECLARE
    table_oid oid := pg_event_trigger_table_rewrite_oid();
    current_hour integer := extract('hour' from current_time);
```

```
pages integer;
max_pages integer := 100;
BEGIN
  IF pg_event_trigger_table_rewrite_oid() = 'public.foo'::regclass
  THEN
    RAISE EXCEPTION 'you're not allowed to rewrite the table %',
      table_oid::regclass;
  END IF;

  SELECT INTO pages relpages FROM pg_class WHERE oid = table_oid;
  IF pages > max_pages
  THEN
    RAISE EXCEPTION 'rewrites only allowed for table with less than % pages',
      max_pages;
  END IF;

  IF current_hour NOT BETWEEN 1 AND 6
  THEN
    RAISE EXCEPTION 'rewrites only allowed between 1am and 6am';
  END IF;
END;
$$;

CREATE EVENT TRIGGER no_rewrite_allowed
  ON table_rewrite
  EXECUTE FUNCTION no_rewrite();
```

41.6. A Database Login Event Trigger Example

The event trigger on the `login` event can be useful for logging user logins, for verifying the connection and assigning roles according to current circumstances, or for session data initialization. It is very important that any event trigger using the `login` event checks whether or not the database is in recovery before performing any writes. Writing to a standby server will make it inaccessible.

The following example demonstrates these options.

```
-- create test tables and roles
CREATE TABLE user_login_log (
  "user" text,
  "session_start" timestamp with time zone
);
CREATE ROLE day_worker;
CREATE ROLE night_worker;

-- the example trigger function
CREATE OR REPLACE FUNCTION init_session()
  RETURNS event_trigger SECURITY DEFINER
  LANGUAGE plpgsql AS
$$
DECLARE
  hour integer = EXTRACT('hour' FROM current_time at time zone 'utc');
  rec boolean;
BEGIN
  -- 1. Forbid logging in between 2AM and 4AM.
  IF hour BETWEEN 2 AND 4 THEN
    RAISE EXCEPTION 'Login forbidden';
  END IF;

  -- The checks below cannot be performed on standby servers so
```

```
-- ensure the database is not in recovery before we perform any
-- operations.
SELECT pg_is_in_recovery() INTO rec;
IF rec THEN
    RETURN;
END IF;

-- 2. Assign some roles. At daytime, grant the day_worker role, else the
-- night_worker role.
IF hour BETWEEN 8 AND 20 THEN
    EXECUTE 'REVOKE night_worker FROM ' || quote_ident(session_user);
    EXECUTE 'GRANT day_worker TO ' || quote_ident(session_user);
ELSE
    EXECUTE 'REVOKE day_worker FROM ' || quote_ident(session_user);
    EXECUTE 'GRANT night_worker TO ' || quote_ident(session_user);
END IF;

-- 3. Initialize user session data
CREATE TEMP TABLE session_storage (x float, y integer);
ALTER TABLE session_storage OWNER TO session_user;

-- 4. Log the connection time
INSERT INTO public.user_login_log VALUES (session_user, current_timestamp);

END;
$$;

-- trigger definition
CREATE EVENT TRIGGER init_session
    ON login
    EXECUTE FUNCTION init_session();
ALTER EVENT TRIGGER init_session ENABLE ALWAYS;
```

Chapter 42. The Rule System

This chapter discusses the rule system in Postgres Pro. Production rule systems are conceptually simple, but there are many subtle points involved in actually using them.

Some other database systems define active database rules, which are usually stored procedures and triggers. In Postgres Pro, these can be implemented using functions and triggers as well.

The rule system (more precisely speaking, the query rewrite rule system) is totally different from stored procedures and triggers. It modifies queries to take rules into consideration, and then passes the modified query to the query planner for planning and execution. It is very powerful, and can be used for many things such as query language procedures, views, and versions. The theoretical foundations and the power of this rule system are also discussed in [ston90b](#) and [ong90](#).

42.1. The Query Tree

To understand how the rule system works it is necessary to know when it is invoked and what its input and results are.

The rule system is located between the parser and the planner. It takes the output of the parser, one query tree, and the user-defined rewrite rules, which are also query trees with some extra information, and creates zero or more query trees as result. So its input and output are always things the parser itself could have produced and thus, anything it sees is basically representable as an SQL statement.

Now what is a query tree? It is an internal representation of an SQL statement where the single parts that it is built from are stored separately. These query trees can be shown in the server log if you set the configuration parameters `debug_print_parse`, `debug_print_rewritten`, or `debug_print_plan`. The rule actions are also stored as query trees, in the system catalog `pg_rewrite`. They are not formatted like the log output, but they contain exactly the same information.

Reading a raw query tree requires some experience. But since SQL representations of query trees are sufficient to understand the rule system, this chapter will not teach how to read them.

When reading the SQL representations of the query trees in this chapter it is necessary to be able to identify the parts the statement is broken into when it is in the query tree structure. The parts of a query tree are

the command type

This is a simple value telling which command (`SELECT`, `INSERT`, `UPDATE`, `DELETE`) produced the query tree.

the range table

The range table is a list of relations that are used in the query. In a `SELECT` statement these are the relations given after the `FROM` key word.

Every range table entry identifies a table or view and tells by which name it is called in the other parts of the query. In the query tree, the range table entries are referenced by number rather than by name, so here it doesn't matter if there are duplicate names as it would in an SQL statement. This can happen after the range tables of rules have been merged in. The examples in this chapter will not have this situation.

the result relation

This is an index into the range table that identifies the relation where the results of the query go.

`SELECT` queries don't have a result relation. (The special case of `SELECT INTO` is mostly identical to `CREATE TABLE` followed by `INSERT ... SELECT`, and is not discussed separately here.)

For `INSERT`, `UPDATE`, and `DELETE` commands, the result relation is the table (or view!) where the changes are to take effect.

the target list

The target list is a list of expressions that define the result of the query. In the case of a `SELECT`, these expressions are the ones that build the final output of the query. They correspond to the expressions between the key words `SELECT` and `FROM`. (* is just an abbreviation for all the column names of a relation. It is expanded by the parser into the individual columns, so the rule system never sees it.)

`DELETE` commands don't need a normal target list because they don't produce any result. Instead, the planner adds a special CTID entry to the empty target list, to allow the executor to find the row to be deleted. (CTID is added when the result relation is an ordinary table. If it is a view, a whole-row variable is added instead, by the rule system, as described in [Section 42.2.4](#).)

For `INSERT` commands, the target list describes the new rows that should go into the result relation. It consists of the expressions in the `VALUES` clause or the ones from the `SELECT` clause in `INSERT ... SELECT`. The first step of the rewrite process adds target list entries for any columns that were not assigned to by the original command but have defaults. Any remaining columns (with neither a given value nor a default) will be filled in by the planner with a constant null expression.

For `UPDATE` commands, the target list describes the new rows that should replace the old ones. In the rule system, it contains just the expressions from the `SET column = expression` part of the command. The planner will handle missing columns by inserting expressions that copy the values from the old row into the new one. Just as for `DELETE`, a CTID or whole-row variable is added so that the executor can identify the old row to be updated.

Every entry in the target list contains an expression that can be a constant value, a variable pointing to a column of one of the relations in the range table, a parameter, or an expression tree made of function calls, constants, variables, operators, etc.

the qualification

The query's qualification is an expression much like one of those contained in the target list entries. The result value of this expression is a Boolean that tells whether the operation (`INSERT`, `UPDATE`, `DELETE`, or `SELECT`) for the final result row should be executed or not. It corresponds to the `WHERE` clause of an SQL statement.

the join tree

The query's join tree shows the structure of the `FROM` clause. For a simple query like `SELECT ... FROM a, b, c`, the join tree is just a list of the `FROM` items, because we are allowed to join them in any order. But when `JOIN` expressions, particularly outer joins, are used, we have to join in the order shown by the joins. In that case, the join tree shows the structure of the `JOIN` expressions. The restrictions associated with particular `JOIN` clauses (from `ON` or `USING` expressions) are stored as qualification expressions attached to those join-tree nodes. It turns out to be convenient to store the top-level `WHERE` expression as a qualification attached to the top-level join-tree item, too. So really the join tree represents both the `FROM` and `WHERE` clauses of a `SELECT`.

the others

The other parts of the query tree like the `ORDER BY` clause aren't of interest here. The rule system substitutes some entries there while applying rules, but that doesn't have much to do with the fundamentals of the rule system.

42.2. Views and the Rule System

Views in Postgres Pro are implemented using the rule system. A view is basically an empty table (having no actual storage) with an `ON SELECT DO INSTEAD` rule. Conventionally, that rule is named `_RETURN`. So a view like

```
CREATE VIEW myview AS SELECT * FROM mytab;
```

is very nearly the same thing as

```
CREATE TABLE myview (same column list as mytab);
CREATE RULE "_RETURN" AS ON SELECT TO myview DO INSTEAD
    SELECT * FROM mytab;
```

although you can't actually write that, because tables are not allowed to have `ON SELECT` rules.

A view can also have other kinds of `DO INSTEAD` rules, allowing `INSERT`, `UPDATE`, or `DELETE` commands to be performed on the view despite its lack of underlying storage. This is discussed further below, in [Section 42.2.4](#).

42.2.1. How `SELECT` Rules Work

Rules `ON SELECT` are applied to all queries as the last step, even if the command given is an `INSERT`, `UPDATE` or `DELETE`. And they have different semantics from rules on the other command types in that they modify the query tree in place instead of creating a new one. So `SELECT` rules are described first.

Currently, there can be only one action in an `ON SELECT` rule, and it must be an unconditional `SELECT` action that is `INSTEAD`. This restriction was required to make rules safe enough to open them for ordinary users, and it restricts `ON SELECT` rules to act like views.

The examples for this chapter are two join views that do some calculations and some more views using them in turn. One of the two first views is customized later by adding rules for `INSERT`, `UPDATE`, and `DELETE` operations so that the final result will be a view that behaves like a real table with some magic functionality. This is not such a simple example to start from and this makes things harder to get into. But it's better to have one example that covers all the points discussed step by step rather than having many different ones that might mix up in mind.

The real tables we need in the first two rule system descriptions are these:

```
CREATE TABLE shoe_data (
    shoenam     text,           -- primary key
    sh_avail    integer,        -- available number of pairs
    slcolor     text,           -- preferred shoelace color
    slminlen    real,           -- minimum shoelace length
    slmaxlen    real,           -- maximum shoelace length
    slunit      text            -- length unit
);

CREATE TABLE shoelace_data (
    sl_name     text,           -- primary key
    sl_avail    integer,        -- available number of pairs
    sl_color    text,           -- shoelace color
    sl_len      real,           -- shoelace length
    sl_unit     text            -- length unit
);

CREATE TABLE unit (
    un_name     text,           -- primary key
    un_fact     real            -- factor to transform to cm
);
```

As you can see, they represent shoe-store data.

The views are created as:

```
CREATE VIEW shoe AS
    SELECT sh.shoenam,
           sh.sh_avail,
           sh.slcolor,
           sh.slminlen,
```

```

        sh.slminlen * un.un_fact AS slminlen_cm,
        sh.slmaxlen,
        sh.slmaxlen * un.un_fact AS slmaxlen_cm,
        sh.slunit
    FROM shoe_data sh, unit un
    WHERE sh.slunit = un.un_name;

```

```

CREATE VIEW shoelace AS
    SELECT s.sl_name,
           s.sl_avail,
           s.sl_color,
           s.sl_len,
           s.sl_unit,
           s.sl_len * u.un_fact AS sl_len_cm
    FROM shoelace_data s, unit u
    WHERE s.sl_unit = u.un_name;

```

```

CREATE VIEW shoe_ready AS
    SELECT rsh.shoename,
           rsh.sh_avail,
           rsl.sl_name,
           rsl.sl_avail,
           least(rsh.sh_avail, rsl.sl_avail) AS total_avail
    FROM shoe rsh, shoelace rsl
    WHERE rsl.sl_color = rsh.slcolor
        AND rsl.sl_len_cm >= rsh.slminlen_cm
        AND rsl.sl_len_cm <= rsh.slmaxlen_cm;

```

The `CREATE VIEW` command for the shoelace view (which is the simplest one we have) will create a relation shoelace and an entry in `pg_rewrite` that tells that there is a rewrite rule that must be applied whenever the relation shoelace is referenced in a query's range table. The rule has no rule qualification (discussed later, with the non-`SELECT` rules, since `SELECT` rules currently cannot have them) and it is `INSTEAD`. Note that rule qualifications are not the same as query qualifications. The action of our rule has a query qualification. The action of the rule is one query tree that is a copy of the `SELECT` statement in the view creation command.

Note

The two extra range table entries for `NEW` and `OLD` that you can see in the `pg_rewrite` entry aren't of interest for `SELECT` rules.

Now we populate `unit`, `shoe_data` and `shoelace_data` and run a simple query on a view:

```

INSERT INTO unit VALUES ('cm', 1.0);
INSERT INTO unit VALUES ('m', 100.0);
INSERT INTO unit VALUES ('inch', 2.54);

INSERT INTO shoe_data VALUES ('sh1', 2, 'black', 70.0, 90.0, 'cm');
INSERT INTO shoe_data VALUES ('sh2', 0, 'black', 30.0, 40.0, 'inch');
INSERT INTO shoe_data VALUES ('sh3', 4, 'brown', 50.0, 65.0, 'cm');
INSERT INTO shoe_data VALUES ('sh4', 3, 'brown', 40.0, 50.0, 'inch');

INSERT INTO shoelace_data VALUES ('sl1', 5, 'black', 80.0, 'cm');
INSERT INTO shoelace_data VALUES ('sl2', 6, 'black', 100.0, 'cm');
INSERT INTO shoelace_data VALUES ('sl3', 0, 'black', 35.0, 'inch');
INSERT INTO shoelace_data VALUES ('sl4', 8, 'black', 40.0, 'inch');
INSERT INTO shoelace_data VALUES ('sl5', 4, 'brown', 1.0, 'm');
INSERT INTO shoelace_data VALUES ('sl6', 0, 'brown', 0.9, 'm');

```

```
INSERT INTO shoelace_data VALUES ('sl7', 7, 'brown', 60 , 'cm');
INSERT INTO shoelace_data VALUES ('sl8', 1, 'brown', 40 , 'inch');
```

```
SELECT * FROM shoelace;
```

sl_name	sl_avail	sl_color	sl_len	sl_unit	sl_len_cm
sl1	5	black	80	cm	80
sl2	6	black	100	cm	100
sl7	7	brown	60	cm	60
sl3	0	black	35	inch	88.9
sl4	8	black	40	inch	101.6
sl8	1	brown	40	inch	101.6
sl5	4	brown	1	m	100
sl6	0	brown	0.9	m	90

(8 rows)

This is the simplest `SELECT` you can do on our views, so we take this opportunity to explain the basics of view rules. The `SELECT * FROM shoelace` was interpreted by the parser and produced the query tree:

```
SELECT shoelace.sl_name, shoelace.sl_avail,
       shoelace.sl_color, shoelace.sl_len,
       shoelace.sl_unit, shoelace.sl_len_cm
FROM shoelace shoelace;
```

and this is given to the rule system. The rule system walks through the range table and checks if there are rules for any relation. When processing the range table entry for `shoelace` (the only one up to now) it finds the `_RETURN` rule with the query tree:

```
SELECT s.sl_name, s.sl_avail,
       s.sl_color, s.sl_len, s.sl_unit,
       s.sl_len * u.un_fact AS sl_len_cm
FROM shoelace old, shoelace new,
     shoelace_data s, unit u
WHERE s.sl_unit = u.un_name;
```

To expand the view, the rewriter simply creates a subquery range-table entry containing the rule's action query tree, and substitutes this range table entry for the original one that referenced the view. The resulting rewritten query tree is almost the same as if you had typed:

```
SELECT shoelace.sl_name, shoelace.sl_avail,
       shoelace.sl_color, shoelace.sl_len,
       shoelace.sl_unit, shoelace.sl_len_cm
FROM (SELECT s.sl_name,
            s.sl_avail,
            s.sl_color,
            s.sl_len,
            s.sl_unit,
            s.sl_len * u.un_fact AS sl_len_cm
      FROM shoelace_data s, unit u
      WHERE s.sl_unit = u.un_name) shoelace;
```

There is one difference however: the subquery's range table has two extra entries `shoelace old` and `shoelace new`. These entries don't participate directly in the query, since they aren't referenced by the subquery's join tree or target list. The rewriter uses them to store the access privilege check information that was originally present in the range-table entry that referenced the view. In this way, the executor will still check that the user has proper privileges to access the view, even though there's no direct use of the view in the rewritten query.

That was the first rule applied. The rule system will continue checking the remaining range-table entries in the top query (in this example there are no more), and it will recursively check the range-table entries

in the added subquery to see if any of them reference views. (But it won't expand `old` or `new` — otherwise we'd have infinite recursion!) In this example, there are no rewrite rules for `shoelace_data` or `unit`, so rewriting is complete and the above is the final result given to the planner.

Now we want to write a query that finds out for which shoes currently in the store we have the matching shoelaces (color and length) and where the total number of exactly matching pairs is greater than or equal to two.

```
SELECT * FROM shoe_ready WHERE total_avail >= 2;
```

shoename	sh_avail	sl_name	sl_avail	total_avail
sh1	2	sl1	5	2
sh3	4	sl7	7	4

(2 rows)

The output of the parser this time is the query tree:

```
SELECT shoe_ready.shoename, shoe_ready.sh_avail,
       shoe_ready.sl_name, shoe_ready.sl_avail,
       shoe_ready.total_avail
FROM shoe_ready shoe_ready
WHERE shoe_ready.total_avail >= 2;
```

The first rule applied will be the one for the `shoe_ready` view and it results in the query tree:

```
SELECT shoe_ready.shoename, shoe_ready.sh_avail,
       shoe_ready.sl_name, shoe_ready.sl_avail,
       shoe_ready.total_avail
FROM (SELECT rsh.shoename,
            rsh.sh_avail,
            rsl.sl_name,
            rsl.sl_avail,
            least(rsh.sh_avail, rsl.sl_avail) AS total_avail
      FROM shoe rsh, shoelace rsl
      WHERE rsl.sl_color = rsh.slcolor
            AND rsl.sl_len_cm >= rsh.slminlen_cm
            AND rsl.sl_len_cm <= rsh.slmaxlen_cm) shoe_ready
WHERE shoe_ready.total_avail >= 2;
```

Similarly, the rules for `shoe` and `shoelace` are substituted into the range table of the subquery, leading to a three-level final query tree:

```
SELECT shoe_ready.shoename, shoe_ready.sh_avail,
       shoe_ready.sl_name, shoe_ready.sl_avail,
       shoe_ready.total_avail
FROM (SELECT rsh.shoename,
            rsh.sh_avail,
            rsl.sl_name,
            rsl.sl_avail,
            least(rsh.sh_avail, rsl.sl_avail) AS total_avail
      FROM (SELECT sh.shoename,
                  sh.sh_avail,
                  sh.slcolor,
                  sh.slminlen,
                  sh.slminlen * un.un_fact AS slminlen_cm,
                  sh.slmaxlen,
                  sh.slmaxlen * un.un_fact AS slmaxlen_cm,
                  sh.slunit
            FROM shoe_data sh, unit un
            WHERE sh.slunit = un.un_name) rsh,
```

```
(SELECT s.sl_name,
       s.sl_avail,
       s.sl_color,
       s.sl_len,
       s.sl_unit,
       s.sl_len * u.un_fact AS sl_len_cm
  FROM shoelace_data s, unit u
 WHERE s.sl_unit = u.un_name) rsl
WHERE rsl.sl_color = rsh.slcolor
  AND rsl.sl_len_cm >= rsh.slminlen_cm
  AND rsl.sl_len_cm <= rsh.slmaxlen_cm) shoe_ready
WHERE shoe_ready.total_avail > 2;
```

This might look inefficient, but the planner will collapse this into a single-level query tree by “pulling up” the subqueries, and then it will plan the joins just as if we’d written them out manually. So collapsing the query tree is an optimization that the rewrite system doesn’t have to concern itself with.

42.2.2. View Rules in Non-SELECT Statements

Two details of the query tree aren’t touched in the description of view rules above. These are the command type and the result relation. In fact, the command type is not needed by view rules, but the result relation may affect the way in which the query rewriter works, because special care needs to be taken if the result relation is a view.

There are only a few differences between a query tree for a `SELECT` and one for any other command. Obviously, they have a different command type and for a command other than a `SELECT`, the result relation points to the range-table entry where the result should go. Everything else is absolutely the same. So having two tables `t1` and `t2` with columns `a` and `b`, the query trees for the two statements:

```
SELECT t2.b FROM t1, t2 WHERE t1.a = t2.a;
```

```
UPDATE t1 SET b = t2.b FROM t2 WHERE t1.a = t2.a;
```

are nearly identical. In particular:

- The range tables contain entries for the tables `t1` and `t2`.
- The target lists contain one variable that points to column `b` of the range table entry for table `t2`.
- The qualification expressions compare the columns `a` of both range-table entries for equality.
- The join trees show a simple join between `t1` and `t2`.

The consequence is, that both query trees result in similar execution plans: They are both joins over the two tables. For the `UPDATE` the missing columns from `t1` are added to the target list by the planner and the final query tree will read as:

```
UPDATE t1 SET a = t1.a, b = t2.b FROM t2 WHERE t1.a = t2.a;
```

and thus the executor run over the join will produce exactly the same result set as:

```
SELECT t1.a, t2.b FROM t1, t2 WHERE t1.a = t2.a;
```

But there is a little problem in `UPDATE`: the part of the executor plan that does the join does not care what the results from the join are meant for. It just produces a result set of rows. The fact that one is a `SELECT` command and the other is an `UPDATE` is handled higher up in the executor, where it knows that this is an `UPDATE`, and it knows that this result should go into table `t1`. But which of the rows that are there has to be replaced by the new row?

To resolve this problem, another entry is added to the target list in `UPDATE` (and also in `DELETE`) statements: the current tuple ID (CTID). This is a system column containing the file block number and position in the block for the row. Knowing the table, the CTID can be used to retrieve the original row of `t1` to be updated. After adding the CTID to the target list, the query actually looks like:

```
SELECT t1.a, t2.b, t1.ctid FROM t1, t2 WHERE t1.a = t2.a;
```

Now another detail of Postgres Pro enters the stage. Old table rows aren't overwritten, and this is why `ROLLBACK` is fast. In an `UPDATE`, the new result row is inserted into the table (after stripping the CTID) and in the row header of the old row, which the CTID pointed to, the `cmax` and `xmax` entries are set to the current command counter and current transaction ID. Thus the old row is hidden, and after the transaction commits the vacuum cleaner can eventually remove the dead row.

Knowing all that, we can simply apply view rules in absolutely the same way to any command. There is no difference.

42.2.3. The Power of Views in Postgres Pro

The above demonstrates how the rule system incorporates view definitions into the original query tree. In the second example, a simple `SELECT` from one view created a final query tree that is a join of 4 tables (`unit` was used twice with different names).

The benefit of implementing views with the rule system is that the planner has all the information about which tables have to be scanned plus the relationships between these tables plus the restrictive qualifications from the views plus the qualifications from the original query in one single query tree. And this is still the situation when the original query is already a join over views. The planner has to decide which is the best path to execute the query, and the more information the planner has, the better this decision can be. And the rule system as implemented in Postgres Pro ensures that this is all information available about the query up to that point.

42.2.4. Updating a View

What happens if a view is named as the target relation for an `INSERT`, `UPDATE`, `DELETE`, or `MERGE`? Doing the substitutions described above would give a query tree in which the result relation points at a subquery range-table entry, which will not work. There are several ways in which Postgres Pro can support the appearance of updating a view, however. In order of user-experienced complexity those are: automatically substitute in the underlying table for the view, execute a user-defined trigger, or rewrite the query per a user-defined rule. These options are discussed below.

If the subquery selects from a single base relation and is simple enough, the rewriter can automatically replace the subquery with the underlying base relation so that the `INSERT`, `UPDATE`, `DELETE`, or `MERGE` is applied to the base relation in the appropriate way. Views that are “simple enough” for this are called *automatically updatable*. For detailed information on the kinds of view that can be automatically updated, see [CREATE VIEW](#).

Alternatively, the operation may be handled by a user-provided `INSTEAD OF` trigger on the view (see [CREATE TRIGGER](#)). Rewriting works slightly differently in this case. For `INSERT`, the rewriter does nothing at all with the view, leaving it as the result relation for the query. For `UPDATE`, `DELETE`, and `MERGE`, it's still necessary to expand the view query to produce the “old” rows that the command will attempt to update, delete, or merge. So the view is expanded as normal, but another unexpanded range-table entry is added to the query to represent the view in its capacity as the result relation.

The problem that now arises is how to identify the rows to be updated in the view. Recall that when the result relation is a table, a special CTID entry is added to the target list to identify the physical locations of the rows to be updated. This does not work if the result relation is a view, because a view does not have any CTID, since its rows do not have actual physical locations. Instead, for an `UPDATE`, `DELETE`, or `MERGE` operation, a special `wholerow` entry is added to the target list, which expands to include all columns from the view. The executor uses this value to supply the “old” row to the `INSTEAD OF` trigger. It is up to the trigger to work out what to update based on the old and new row values.

Another possibility is for the user to define `INSTEAD` rules that specify substitute actions for `INSERT`, `UPDATE`, and `DELETE` commands on a view. These rules will rewrite the command, typically into a command that updates one or more tables, rather than views. That is the topic of [Section 42.4](#). Note that this will not work with `MERGE`, which currently does not support rules on the target relation other than `SELECT` rules.

Note that rules are evaluated first, rewriting the original query before it is planned and executed. Therefore, if a view has `INSTEAD OF` triggers as well as rules on `INSERT`, `UPDATE`, or `DELETE`, then the rules will be evaluated first, and depending on the result, the triggers may not be used at all.

Automatic rewriting of an `INSERT`, `UPDATE`, `DELETE`, or `MERGE` query on a simple view is always tried last. Therefore, if a view has rules or triggers, they will override the default behavior of automatically updatable views.

If there are no `INSTEAD` rules or `INSTEAD OF` triggers for the view, and the rewriter cannot automatically rewrite the query as an update on the underlying base relation, an error will be thrown because the executor cannot update a view as such.

42.3. Materialized Views

Materialized views in Postgres Pro use the rule system like views do, but persist the results in a table-like form. The main differences between:

```
CREATE MATERIALIZED VIEW mymatview AS SELECT * FROM mytab;
```

and:

```
CREATE TABLE mymatview AS SELECT * FROM mytab;
```

are that the materialized view cannot subsequently be directly updated and that the query used to create the materialized view is stored in exactly the same way that a view's query is stored, so that fresh data can be generated for the materialized view with:

```
REFRESH MATERIALIZED VIEW mymatview;
```

The information about a materialized view in the Postgres Pro system catalogs is exactly the same as it is for a table or view. So for the parser, a materialized view is a relation, just like a table or a view. When a materialized view is referenced in a query, the data is returned directly from the materialized view, like from a table; the rule is only used for populating the materialized view.

While access to the data stored in a materialized view is often much faster than accessing the underlying tables directly or through a view, the data is not always current; yet sometimes current data is not needed. Consider a table which records sales:

```
CREATE TABLE invoice (  
    invoice_no    integer          PRIMARY KEY,  
    seller_no     integer,         -- ID of salesperson  
    invoice_date  date,            -- date of sale  
    invoice_amt   numeric(13,2)    -- amount of sale  
);
```

If people want to be able to quickly graph historical sales data, they might want to summarize, and they may not care about the incomplete data for the current date:

```
CREATE MATERIALIZED VIEW sales_summary AS  
SELECT  
    seller_no,  
    invoice_date,  
    sum(invoice_amt)::numeric(13,2) as sales_amt  
FROM invoice  
WHERE invoice_date < CURRENT_DATE  
GROUP BY  
    seller_no,  
    invoice_date;
```

```
CREATE UNIQUE INDEX sales_summary_seller  
ON sales_summary (seller_no, invoice_date);
```

This materialized view might be useful for displaying a graph in the dashboard created for salespeople. A job could be scheduled to update the statistics each night using this SQL statement:

```
REFRESH MATERIALIZED VIEW sales_summary;
```

Another use for a materialized view is to allow faster access to data brought across from a remote system through a foreign data wrapper. A simple example using `file_fdw` is below, with timings, but since this is using cache on the local system the performance difference compared to access to a remote system would usually be greater than shown here. Notice we are also exploiting the ability to put an index on the materialized view, whereas `file_fdw` does not support indexes; this advantage might not apply for other sorts of foreign data access.

Setup:

```
CREATE EXTENSION file_fdw;
CREATE SERVER local_file FOREIGN DATA WRAPPER file_fdw;
CREATE FOREIGN TABLE words (word text NOT NULL)
    SERVER local_file
    OPTIONS (filename '/usr/share/dict/words');
CREATE MATERIALIZED VIEW wrd AS SELECT * FROM words;
CREATE UNIQUE INDEX wrd_word ON wrd (word);
CREATE EXTENSION pg_trgm;
CREATE INDEX wrd_trgm ON wrd USING gist (word gist_trgm_ops);
VACUUM ANALYZE wrd;
```

Now let's spell-check a word. Using `file_fdw` directly:

```
SELECT count(*) FROM words WHERE word = 'caterpiler';
```

```
count
-----
      0
(1 row)
```

With `EXPLAIN ANALYZE`, we see:

```
Aggregate  (cost=21763.99..21764.00 rows=1 width=0) (actual time=188.180..188.181
rows=1 loops=1)
  -> Foreign Scan on words  (cost=0.00..21761.41 rows=1032 width=0) (actual
time=188.177..188.177 rows=0 loops=1)
    Filter: (word = 'caterpiler'::text)
    Rows Removed by Filter: 479829
    Foreign File: /usr/share/dict/words
    Foreign File Size: 4953699
Planning time: 0.118 ms
Execution time: 188.273 ms
```

If the materialized view is used instead, the query is much faster:

```
Aggregate  (cost=4.44..4.45 rows=1 width=0) (actual time=0.042..0.042 rows=1 loops=1)
  -> Index Only Scan using wrd_word on wrd  (cost=0.42..4.44 rows=1 width=0) (actual
time=0.039..0.039 rows=0 loops=1)
    Index Cond: (word = 'caterpiler'::text)
    Heap Fetches: 0
Planning time: 0.164 ms
Execution time: 0.117 ms
```

Either way, the word is spelled wrong, so let's look for what we might have wanted. Again using `file_fdw` and `pg_trgm`:

```
SELECT word FROM words ORDER BY word <-> 'caterpiler' LIMIT 10;
```

```
word
-----
cater
```

```
caterpillar
Caterpillar
caterpillars
caterpillar's
Caterpillar's
caterer
caterer's
caters
catered
(10 rows)

Limit (cost=11583.61..11583.64 rows=10 width=32) (actual time=1431.591..1431.594
rows=10 loops=1)
  -> Sort (cost=11583.61..11804.76 rows=88459 width=32) (actual
time=1431.589..1431.591 rows=10 loops=1)
    Sort Key: ((word <-> 'caterpiler'::text))
    Sort Method: top-N heapsort  Memory: 25kB
  -> Foreign Scan on words (cost=0.00..9672.05 rows=88459 width=32) (actual
time=0.057..1286.455 rows=479829 loops=1)
    Foreign File: /usr/share/dict/words
    Foreign File Size: 4953699
Planning time: 0.128 ms
Execution time: 1431.679 ms
```

Using the materialized view:

```
Limit (cost=0.29..1.06 rows=10 width=10) (actual time=187.222..188.257 rows=10
loops=1)
  -> Index Scan using wrd_trgm on wrd (cost=0.29..37020.87 rows=479829 width=10)
(actual time=187.219..188.252 rows=10 loops=1)
    Order By: (word <-> 'caterpiler'::text)
Planning time: 0.196 ms
Execution time: 198.640 ms
```

If you can tolerate periodic update of the remote data to the local database, the performance benefit can be substantial.

42.4. Rules on INSERT, UPDATE, and DELETE

Rules that are defined on `INSERT`, `UPDATE`, and `DELETE` are significantly different from the view rules described in the previous sections. First, their `CREATE RULE` command allows more:

- They are allowed to have no action.
- They can have multiple actions.
- They can be `INSTEAD` or `ALSO` (the default).
- The pseudorelations `NEW` and `OLD` become useful.
- They can have rule qualifications.

Second, they don't modify the query tree in place. Instead they create zero or more new query trees and can throw away the original one.

Caution

In many cases, tasks that could be performed by rules on `INSERT/UPDATE/DELETE` are better done with triggers. Triggers are notationally a bit more complicated, but their semantics are much simpler to understand. Rules tend to have surprising results when the original query contains volatile functions: volatile functions may get executed more times than expected in the process of carrying out the rules.

Also, there are some cases that are not supported by these types of rules at all, notably including `WITH` clauses in the original query and multiple-assignment sub-`SELECT`s in the `SET` list of `UPDATE` queries. This is because copying these constructs into a rule query would result in multiple evaluations of the sub-query, contrary to the express intent of the query's author.

42.4.1. How Update Rules Work

Keep the syntax:

```
CREATE [ OR REPLACE ] RULE name AS ON event
  TO table [ WHERE condition ]
  DO [ ALSO | INSTEAD ] { NOTHING | command | ( command ; command ... ) }
```

in mind. In the following, *update rules* means rules that are defined on `INSERT`, `UPDATE`, or `DELETE`.

Update rules get applied by the rule system when the result relation and the command type of a query tree are equal to the object and event given in the `CREATE RULE` command. For update rules, the rule system creates a list of query trees. Initially the query-tree list is empty. There can be zero (`NOTHING` key word), one, or multiple actions. To simplify, we will look at a rule with one action. This rule can have a qualification or not and it can be `INSTEAD` or `ALSO` (the default).

What is a rule qualification? It is a restriction that tells when the actions of the rule should be done and when not. This qualification can only reference the pseudorelations `NEW` and/or `OLD`, which basically represent the relation that was given as object (but with a special meaning).

So we have three cases that produce the following query trees for a one-action rule.

No qualification, with either `ALSO` or `INSTEAD`

the query tree from the rule action with the original query tree's qualification added

Qualification given and `ALSO`

the query tree from the rule action with the rule qualification and the original query tree's qualification added

Qualification given and `INSTEAD`

the query tree from the rule action with the rule qualification and the original query tree's qualification; and the original query tree with the negated rule qualification added

Finally, if the rule is `ALSO`, the unchanged original query tree is added to the list. Since only qualified `INSTEAD` rules already add the original query tree, we end up with either one or two output query trees for a rule with one action.

For `ON INSERT` rules, the original query (if not suppressed by `INSTEAD`) is done before any actions added by rules. This allows the actions to see the inserted row(s). But for `ON UPDATE` and `ON DELETE` rules, the original query is done after the actions added by rules. This ensures that the actions can see the to-be-updated or to-be-deleted rows; otherwise, the actions might do nothing because they find no rows matching their qualifications.

The query trees generated from rule actions are thrown into the rewrite system again, and maybe more rules get applied resulting in additional or fewer query trees. So a rule's actions must have either a different command type or a different result relation than the rule itself is on, otherwise this recursive process will end up in an infinite loop. (Recursive expansion of a rule will be detected and reported as an error.)

The query trees found in the actions of the `pg_rewrite` system catalog are only templates. Since they can reference the range-table entries for `NEW` and `OLD`, some substitutions have to be made before they can be used. For any reference to `NEW`, the target list of the original query is searched for a corresponding entry. If found, that entry's expression replaces the reference. Otherwise, `NEW` means the same as `OLD`.

(for an `UPDATE`) or is replaced by a null value (for an `INSERT`). Any reference to `OLD` is replaced by a reference to the range-table entry that is the result relation.

After the system is done applying update rules, it applies view rules to the produced query tree(s). Views cannot insert new update actions so there is no need to apply update rules to the output of view rewriting.

42.4.1.1. A First Rule Step by Step

Say we want to trace changes to the `sl_avail` column in the `shoelace_data` relation. So we set up a log table and a rule that conditionally writes a log entry when an `UPDATE` is performed on `shoelace_data`.

```
CREATE TABLE shoelace_log (  
    sl_name      text,          -- shoelace changed  
    sl_avail     integer,       -- new available value  
    log_who      text,          -- who did it  
    log_when     timestamp      -- when  
);  
  
CREATE RULE log_shoelace AS ON UPDATE TO shoelace_data  
    WHERE NEW.sl_avail <> OLD.sl_avail  
    DO INSERT INTO shoelace_log VALUES (  
        NEW.sl_name,  
        NEW.sl_avail,  
        current_user,  
        current_timestamp  
    );
```

Now someone does:

```
UPDATE shoelace_data SET sl_avail = 6 WHERE sl_name = 'sl7';
```

and we look at the log table:

```
SELECT * FROM shoelace_log;
```

sl_name	sl_avail	log_who	log_when
sl7	6	Al	Tue Oct 20 16:14:45 1998 MET DST

(1 row)

That's what we expected. What happened in the background is the following. The parser created the query tree:

```
UPDATE shoelace_data SET sl_avail = 6  
    FROM shoelace_data shoelace_data  
    WHERE shoelace_data.sl_name = 'sl7';
```

There is a rule `log_shoelace` that is `ON UPDATE` with the rule qualification expression:

```
NEW.sl_avail <> OLD.sl_avail
```

and the action:

```
INSERT INTO shoelace_log VALUES (  
    new.sl_name, new.sl_avail,  
    current_user, current_timestamp )  
FROM shoelace_data new, shoelace_data old;
```

(This looks a little strange since you cannot normally write `INSERT ... VALUES ... FROM`. The `FROM` clause here is just to indicate that there are range-table entries in the query tree for `new` and `old`. These are needed so that they can be referenced by variables in the `INSERT` command's query tree.)

The rule is a qualified `ALSO` rule, so the rule system has to return two query trees: the modified rule action and the original query tree. In step 1, the range table of the original query is incorporated into the rule's action query tree. This results in:


```
INSERT INTO shoelace_log VALUES (
    new.sl_name, new.sl_avail,
    current_user, current_timestamp )
FROM shoelace_data new, shoelace_data old,
    shoelace_data shoelace_data;
```

In step 2, the rule qualification is added to it, so the result set is restricted to rows where `sl_avail` changes:

```
INSERT INTO shoelace_log VALUES (
    new.sl_name, new.sl_avail,
    current_user, current_timestamp )
FROM shoelace_data new, shoelace_data old,
    shoelace_data shoelace_data
WHERE new.sl_avail <> old.sl_avail;
```

(This looks even stranger, since `INSERT ... VALUES` doesn't have a `WHERE` clause either, but the planner and executor will have no difficulty with it. They need to support this same functionality anyway for `INSERT ... SELECT`.)

In step 3, the original query tree's qualification is added, restricting the result set further to only the rows that would have been touched by the original query:

```
INSERT INTO shoelace_log VALUES (
    new.sl_name, new.sl_avail,
    current_user, current_timestamp )
FROM shoelace_data new, shoelace_data old,
    shoelace_data shoelace_data
WHERE new.sl_avail <> old.sl_avail
    AND shoelace_data.sl_name = 'sl7';
```

Step 4 replaces references to `NEW` by the target list entries from the original query tree or by the matching variable references from the result relation:

```
INSERT INTO shoelace_log VALUES (
    shoelace_data.sl_name, 6,
    current_user, current_timestamp )
FROM shoelace_data new, shoelace_data old,
    shoelace_data shoelace_data
WHERE 6 <> old.sl_avail
    AND shoelace_data.sl_name = 'sl7';
```

Step 5 changes `OLD` references into result relation references:

```
INSERT INTO shoelace_log VALUES (
    shoelace_data.sl_name, 6,
    current_user, current_timestamp )
FROM shoelace_data new, shoelace_data old,
    shoelace_data shoelace_data
WHERE 6 <> shoelace_data.sl_avail
    AND shoelace_data.sl_name = 'sl7';
```

That's it. Since the rule is `ALSO`, we also output the original query tree. In short, the output from the rule system is a list of two query trees that correspond to these statements:

```
INSERT INTO shoelace_log VALUES (
    shoelace_data.sl_name, 6,
    current_user, current_timestamp )
FROM shoelace_data
WHERE 6 <> shoelace_data.sl_avail
    AND shoelace_data.sl_name = 'sl7';
```

```
UPDATE shoelace_data SET sl_avail = 6
WHERE sl_name = 'sl7';
```

These are executed in this order, and that is exactly what the rule was meant to do.

The substitutions and the added qualifications ensure that, if the original query would be, say:

```
UPDATE shoelace_data SET sl_color = 'green'
WHERE sl_name = 'sl7';
```

no log entry would get written. In that case, the original query tree does not contain a target list entry for `sl_avail`, so `NEW.sl_avail` will get replaced by `shoelace_data.sl_avail`. Thus, the extra command generated by the rule is:

```
INSERT INTO shoelace_log VALUES (
    shoelace_data.sl_name, shoelace_data.sl_avail,
    current_user, current_timestamp )
FROM shoelace_data
WHERE shoelace_data.sl_avail <> shoelace_data.sl_avail
AND shoelace_data.sl_name = 'sl7';
```

and that qualification will never be true.

It will also work if the original query modifies multiple rows. So if someone issued the command:

```
UPDATE shoelace_data SET sl_avail = 0
WHERE sl_color = 'black';
```

four rows in fact get updated (`sl1`, `sl2`, `sl3`, and `sl4`). But `sl3` already has `sl_avail = 0`. In this case, the original query tree's qualification is different and that results in the extra query tree:

```
INSERT INTO shoelace_log
SELECT shoelace_data.sl_name, 0,
       current_user, current_timestamp
FROM shoelace_data
WHERE 0 <> shoelace_data.sl_avail
AND shoelace_data.sl_color = 'black';
```

being generated by the rule. This query tree will surely insert three new log entries. And that's absolutely correct.

Here we can see why it is important that the original query tree is executed last. If the `UPDATE` had been executed first, all the rows would have already been set to zero, so the logging `INSERT` would not find any row where `0 <> shoelace_data.sl_avail`.

42.4.2. Cooperation with Views

A simple way to protect view relations from the mentioned possibility that someone can try to run `INSERT`, `UPDATE`, or `DELETE` on them is to let those query trees get thrown away. So we could create the rules:

```
CREATE RULE shoe_ins_protect AS ON INSERT TO shoe
DO INSTEAD NOTHING;
CREATE RULE shoe_upd_protect AS ON UPDATE TO shoe
DO INSTEAD NOTHING;
CREATE RULE shoe_del_protect AS ON DELETE TO shoe
DO INSTEAD NOTHING;
```

If someone now tries to do any of these operations on the view relation `shoe`, the rule system will apply these rules. Since the rules have no actions and are `INSTEAD`, the resulting list of query trees will be empty and the whole query will become nothing because there is nothing left to be optimized or executed after the rule system is done with it.

A more sophisticated way to use the rule system is to create rules that rewrite the query tree into one that does the right operation on the real tables. To do that on the `shoelace` view, we create the following rules:

```
CREATE RULE shoelace_ins AS ON INSERT TO shoelace
DO INSTEAD
INSERT INTO shoelace_data VALUES (
    NEW.sl_name,
    NEW.sl_avail,
    NEW.sl_color,
    NEW.sl_len,
    NEW.sl_unit
);
```

```
CREATE RULE shoelace_upd AS ON UPDATE TO shoelace
DO INSTEAD
UPDATE shoelace_data
SET sl_name = NEW.sl_name,
    sl_avail = NEW.sl_avail,
    sl_color = NEW.sl_color,
    sl_len = NEW.sl_len,
    sl_unit = NEW.sl_unit
WHERE sl_name = OLD.sl_name;
```

```
CREATE RULE shoelace_del AS ON DELETE TO shoelace
DO INSTEAD
DELETE FROM shoelace_data
WHERE sl_name = OLD.sl_name;
```

If you want to support RETURNING queries on the view, you need to make the rules include RETURNING clauses that compute the view rows. This is usually pretty trivial for views on a single table, but it's a bit tedious for join views such as shoelace. An example for the insert case is:

```
CREATE RULE shoelace_ins AS ON INSERT TO shoelace
DO INSTEAD
INSERT INTO shoelace_data VALUES (
    NEW.sl_name,
    NEW.sl_avail,
    NEW.sl_color,
    NEW.sl_len,
    NEW.sl_unit
)
RETURNING
    shoelace_data.*,
    (SELECT shoelace_data.sl_len * u.un_fact
     FROM unit u WHERE shoelace_data.sl_unit = u.un_name);
```

Note that this one rule supports both INSERT and INSERT RETURNING queries on the view — the RETURNING clause is simply ignored for INSERT.

Now assume that once in a while, a pack of shoelaces arrives at the shop and a big parts list along with it. But you don't want to manually update the shoelace view every time. Instead we set up two little tables: one where you can insert the items from the part list, and one with a special trick. The creation commands for these are:

```
CREATE TABLE shoelace_arrive (
    arr_name    text,
    arr_quant   integer
);

CREATE TABLE shoelace_ok (
    ok_name     text,
    ok_quant    integer
);
```

```
CREATE RULE shoelace_ok_ins AS ON INSERT TO shoelace_ok
DO INSTEAD
UPDATE shoelace
SET sl_avail = sl_avail + NEW.ok_quant
WHERE sl_name = NEW.ok_name;
```

Now you can fill the table shoelace_arrive with the data from the parts list:

```
SELECT * FROM shoelace_arrive;
```

arr_name	arr_quant
sl3	10
sl6	20
sl8	20

(3 rows)

Take a quick look at the current data:

```
SELECT * FROM shoelace;
```

sl_name	sl_avail	sl_color	sl_len	sl_unit	sl_len_cm
sl1	5	black	80	cm	80
sl2	6	black	100	cm	100
sl7	6	brown	60	cm	60
sl3	0	black	35	inch	88.9
sl4	8	black	40	inch	101.6
sl8	1	brown	40	inch	101.6
sl5	4	brown	1	m	100
sl6	0	brown	0.9	m	90

(8 rows)

Now move the arrived shoelaces in:

```
INSERT INTO shoelace_ok SELECT * FROM shoelace_arrive;
```

and check the results:

```
SELECT * FROM shoelace ORDER BY sl_name;
```

sl_name	sl_avail	sl_color	sl_len	sl_unit	sl_len_cm
sl1	5	black	80	cm	80
sl2	6	black	100	cm	100
sl7	6	brown	60	cm	60
sl4	8	black	40	inch	101.6
sl3	10	black	35	inch	88.9
sl8	21	brown	40	inch	101.6
sl5	4	brown	1	m	100
sl6	20	brown	0.9	m	90

(8 rows)

```
SELECT * FROM shoelace_log;
```

sl_name	sl_avail	log_who	log_when
sl7	6	Al	Tue Oct 20 19:14:45 1998 MET DST
sl3	10	Al	Tue Oct 20 19:25:16 1998 MET DST
sl6	20	Al	Tue Oct 20 19:25:16 1998 MET DST
sl8	21	Al	Tue Oct 20 19:25:16 1998 MET DST

(4 rows)

It's a long way from the one `INSERT ... SELECT` to these results. And the description of the query-tree transformation will be the last in this chapter. First, there is the parser's output:

```
INSERT INTO shoelace_ok
SELECT shoelace_arrive.arr_name, shoelace_arrive.arr_quant
FROM shoelace_arrive shoelace_arrive, shoelace_ok shoelace_ok;
```

Now the first rule `shoelace_ok_ins` is applied and turns this into:

```
UPDATE shoelace
SET sl_avail = shoelace.sl_avail + shoelace_arrive.arr_quant
FROM shoelace_arrive shoelace_arrive, shoelace_ok shoelace_ok,
     shoelace_ok old, shoelace_ok new,
     shoelace shoelace
WHERE shoelace.sl_name = shoelace_arrive.arr_name;
```

and throws away the original `INSERT` on `shoelace_ok`. This rewritten query is passed to the rule system again, and the second applied rule `shoelace_upd` produces:

```
UPDATE shoelace_data
SET sl_name = shoelace.sl_name,
    sl_avail = shoelace.sl_avail + shoelace_arrive.arr_quant,
    sl_color = shoelace.sl_color,
    sl_len = shoelace.sl_len,
    sl_unit = shoelace.sl_unit
FROM shoelace_arrive shoelace_arrive, shoelace_ok shoelace_ok,
     shoelace_ok old, shoelace_ok new,
     shoelace shoelace, shoelace old,
     shoelace new, shoelace_data shoelace_data
WHERE shoelace.sl_name = shoelace_arrive.arr_name
     AND shoelace_data.sl_name = shoelace.sl_name;
```

Again it's an `INSTEAD` rule and the previous query tree is trashed. Note that this query still uses the view `shoelace`. But the rule system isn't finished with this step, so it continues and applies the `_RETURN` rule on it, and we get:

```
UPDATE shoelace_data
SET sl_name = s.sl_name,
    sl_avail = s.sl_avail + shoelace_arrive.arr_quant,
    sl_color = s.sl_color,
    sl_len = s.sl_len,
    sl_unit = s.sl_unit
FROM shoelace_arrive shoelace_arrive, shoelace_ok shoelace_ok,
     shoelace_ok old, shoelace_ok new,
     shoelace shoelace, shoelace old,
     shoelace new, shoelace_data shoelace_data,
     shoelace old, shoelace new,
     shoelace_data s, unit u
WHERE s.sl_name = shoelace_arrive.arr_name
     AND shoelace_data.sl_name = s.sl_name;
```

Finally, the rule `log_shoelace` gets applied, producing the extra query tree:

```
INSERT INTO shoelace_log
SELECT s.sl_name,
       s.sl_avail + shoelace_arrive.arr_quant,
       current_user,
       current_timestamp
FROM shoelace_arrive shoelace_arrive, shoelace_ok shoelace_ok,
     shoelace_ok old, shoelace_ok new,
     shoelace shoelace, shoelace old,
```

```

        shoelace new, shoelace_data shoelace_data,
        shoelace old, shoelace new,
        shoelace_data s, unit u,
        shoelace_data old, shoelace_data new
        shoelace_log shoelace_log
WHERE s.sl_name = shoelace_arrive.arr_name
      AND shoelace_data.sl_name = s.sl_name
      AND (s.sl_avail + shoelace_arrive.arr_quant) <> s.sl_avail;

```

After that the rule system runs out of rules and returns the generated query trees.

So we end up with two final query trees that are equivalent to the SQL statements:

```

INSERT INTO shoelace_log
SELECT s.sl_name,
       s.sl_avail + shoelace_arrive.arr_quant,
       current_user,
       current_timestamp
FROM shoelace_arrive shoelace_arrive, shoelace_data shoelace_data,
     shoelace_data s
WHERE s.sl_name = shoelace_arrive.arr_name
      AND shoelace_data.sl_name = s.sl_name
      AND s.sl_avail + shoelace_arrive.arr_quant <> s.sl_avail;

UPDATE shoelace_data
SET sl_avail = shoelace_data.sl_avail + shoelace_arrive.arr_quant
FROM shoelace_arrive shoelace_arrive,
     shoelace_data shoelace_data,
     shoelace_data s
WHERE s.sl_name = shoelace_arrive.sl_name
      AND shoelace_data.sl_name = s.sl_name;

```

The result is that data coming from one relation inserted into another, changed into updates on a third, changed into updating a fourth plus logging that final update in a fifth gets reduced into two queries.

There is a little detail that's a bit ugly. Looking at the two queries, it turns out that the `shoelace_data` relation appears twice in the range table where it could definitely be reduced to one. The planner does not handle it and so the execution plan for the rule systems output of the `INSERT` will be

```

Nested Loop
-> Merge Join
    -> Seq Scan
        -> Sort
            -> Seq Scan on s
    -> Seq Scan
        -> Sort
            -> Seq Scan on shoelace_arrive
-> Seq Scan on shoelace_data

```

while omitting the extra range table entry would result in a

```

Merge Join
-> Seq Scan
    -> Sort
        -> Seq Scan on s
-> Seq Scan
    -> Sort
        -> Seq Scan on shoelace_arrive

```

which produces exactly the same entries in the log table. Thus, the rule system caused one extra scan on the table `shoelace_data` that is absolutely not necessary. And the same redundant scan is done once more in the `UPDATE`. But it was a really hard job to make that all possible at all.

Now we make a final demonstration of the Postgres Pro rule system and its power. Say you add some shoelaces with extraordinary colors to your database:

```
INSERT INTO shoelace VALUES ('sl9', 0, 'pink', 35.0, 'inch', 0.0);
INSERT INTO shoelace VALUES ('sl10', 1000, 'magenta', 40.0, 'inch', 0.0);
```

We would like to make a view to check which shoelace entries do not fit any shoe in color. The view for this is:

```
CREATE VIEW shoelace_mismatch AS
  SELECT * FROM shoelace WHERE NOT EXISTS
    (SELECT shoenam FROM shoe WHERE slcolor = sl_color);
```

Its output is:

```
SELECT * FROM shoelace_mismatch;
```

sl_name	sl_avail	sl_color	sl_len	sl_unit	sl_len_cm
sl9	0	pink	35	inch	88.9
sl10	1000	magenta	40	inch	101.6

Now we want to set it up so that mismatching shoelaces that are not in stock are deleted from the database. To make it a little harder for Postgres Pro, we don't delete it directly. Instead we create one more view:

```
CREATE VIEW shoelace_can_delete AS
  SELECT * FROM shoelace_mismatch WHERE sl_avail = 0;
```

and do it this way:

```
DELETE FROM shoelace WHERE EXISTS
  (SELECT * FROM shoelace_can_delete
    WHERE sl_name = shoelace.sl_name);
```

The results are:

```
SELECT * FROM shoelace;
```

sl_name	sl_avail	sl_color	sl_len	sl_unit	sl_len_cm
sl1	5	black	80	cm	80
sl2	6	black	100	cm	100
sl7	6	brown	60	cm	60
sl4	8	black	40	inch	101.6
sl3	10	black	35	inch	88.9
sl8	21	brown	40	inch	101.6
sl10	1000	magenta	40	inch	101.6
sl5	4	brown	1	m	100
sl6	20	brown	0.9	m	90

(9 rows)

A **DELETE** on a view, with a subquery qualification that in total uses 4 nesting/joined views, where one of them itself has a subquery qualification containing a view and where calculated view columns are used, gets rewritten into one single query tree that deletes the requested data from a real table.

There are probably only a few situations out in the real world where such a construct is necessary. But it makes you feel comfortable that it works.

42.5. Rules and Privileges

Due to rewriting of queries by the Postgres Pro rule system, other tables/views than those used in the original query get accessed. When update rules are used, this can include write access to tables.

Rewrite rules don't have a separate owner. The owner of a relation (table or view) is automatically the owner of the rewrite rules that are defined for it. The Postgres Pro rule system changes the behavior of the default access control system. With the exception of `SELECT` rules associated with security invoker views (see [CREATE VIEW](#)), all relations that are used due to rules get checked against the privileges of the rule owner, not the user invoking the rule. This means that, except for security invoker views, users only need the required privileges for the tables/views that are explicitly named in their queries.

For example: A user has a list of phone numbers where some of them are private, the others are of interest for the assistant of the office. The user can construct the following:

```
CREATE TABLE phone_data (person text, phone text, private boolean);
CREATE VIEW phone_number AS
    SELECT person, CASE WHEN NOT private THEN phone END AS phone
    FROM phone_data;
GRANT SELECT ON phone_number TO assistant;
```

Nobody except that user (and the database superusers) can access the `phone_data` table. But because of the `GRANT`, the assistant can run a `SELECT` on the `phone_number` view. The rule system will rewrite the `SELECT` from `phone_number` into a `SELECT` from `phone_data`. Since the user is the owner of `phone_number` and therefore the owner of the rule, the read access to `phone_data` is now checked against the user's privileges and the query is permitted. The check for accessing `phone_number` is also performed, but this is done against the invoking user, so nobody but the user and the assistant can use it.

The privileges are checked rule by rule. So the assistant is for now the only one who can see the public phone numbers. But the assistant can set up another view and grant access to that to the public. Then, anyone can see the `phone_number` data through the assistant's view. What the assistant cannot do is to create a view that directly accesses `phone_data`. (Actually the assistant can, but it will not work since every access will be denied during the permission checks.) And as soon as the user notices that the assistant opened their `phone_number` view, the user can revoke the assistant's access. Immediately, any access to the assistant's view would fail.

One might think that this rule-by-rule checking is a security hole, but in fact it isn't. But if it did not work this way, the assistant could set up a table with the same columns as `phone_number` and copy the data to there once per day. Then it's the assistant's own data and the assistant can grant access to everyone they want. A `GRANT` command means, "I trust you". If someone you trust does the thing above, it's time to think it over and then use `REVOKE`.

Note that while views can be used to hide the contents of certain columns using the technique shown above, they cannot be used to reliably conceal the data in unseen rows unless the `security_barrier` flag has been set. For example, the following view is insecure:

```
CREATE VIEW phone_number AS
    SELECT person, phone FROM phone_data WHERE phone NOT LIKE '412%';
```

This view might seem secure, since the rule system will rewrite any `SELECT` from `phone_number` into a `SELECT` from `phone_data` and add the qualification that only entries where `phone` does not begin with 412 are wanted. But if the user can create their own functions, it is not difficult to convince the planner to execute the user-defined function prior to the `NOT LIKE` expression. For example:

```
CREATE FUNCTION tricky(text, text) RETURNS bool AS $$
BEGIN
    RAISE NOTICE '% => %', $1, $2;
    RETURN true;
END;
$$ LANGUAGE plpgsql COST 0.000000000000000000000001;

SELECT * FROM phone_number WHERE tricky(person, phone);
```

Every person and phone number in the `phone_data` table will be printed as a `NOTICE`, because the planner will choose to execute the inexpensive `tricky` function before the more expensive `NOT LIKE`. Even if the user is prevented from defining new functions, built-in functions can be used in similar attacks. (For example, most casting functions include their input values in the error messages they produce.)

Similar considerations apply to update rules. In the examples of the previous section, the owner of the tables in the example database could grant the privileges `SELECT`, `INSERT`, `UPDATE`, and `DELETE` on the `shoelace` view to someone else, but only `SELECT` on `shoelace_log`. The rule action to write log entries will still be executed successfully, and that other user could see the log entries. But they could not create fake entries, nor could they manipulate or remove existing ones. In this case, there is no possibility of subverting the rules by convincing the planner to alter the order of operations, because the only rule which references `shoelace_log` is an unqualified `INSERT`. This might not be true in more complex scenarios.

When it is necessary for a view to provide row-level security, the `security_barrier` attribute should be applied to the view. This prevents maliciously-chosen functions and operators from being passed values from rows until after the view has done its work. For example, if the view shown above had been created like this, it would be secure:

```
CREATE VIEW phone_number WITH (security_barrier) AS
  SELECT person, phone FROM phone_data WHERE phone NOT LIKE '412%';
```

Views created with the `security_barrier` may perform far worse than views created without this option. In general, there is no way to avoid this: the fastest possible plan must be rejected if it may compromise security. For this reason, this option is not enabled by default.

The query planner has more flexibility when dealing with functions that have no side effects. Such functions are referred to as `LEAKPROOF`, and include many simple, commonly used operators, such as many equality operators. The query planner can safely allow such functions to be evaluated at any point in the query execution process, since invoking them on rows invisible to the user will not leak any information about the unseen rows. Further, functions which do not take arguments or which are not passed any arguments from the security barrier view do not have to be marked as `LEAKPROOF` to be pushed down, as they never receive data from the view. In contrast, a function that might throw an error depending on the values received as arguments (such as one that throws an error in the event of overflow or division by zero) is not leak-proof, and could provide significant information about the unseen rows if applied before the security view's row filters.

It is important to understand that even a view created with the `security_barrier` option is intended to be secure only in the limited sense that the contents of the invisible tuples will not be passed to possibly-insecure functions. The user may well have other means of making inferences about the unseen data; for example, they can see the query plan using `EXPLAIN`, or measure the run time of queries against the view. A malicious attacker might be able to infer something about the amount of unseen data, or even gain some information about the data distribution or most common values (since these things may affect the run time of the plan; or even, since they are also reflected in the optimizer statistics, the choice of plan). If these types of "covert channel" attacks are of concern, it is probably unwise to grant any access to the data at all.

42.6. Rules and Command Status

The Postgres Pro server returns a command status string, such as `INSERT 149592 1`, for each command it receives. This is simple enough when there are no rules involved, but what happens when the query is rewritten by rules?

Rules affect the command status as follows:

- If there is no unconditional `INSTEAD` rule for the query, then the originally given query will be executed, and its command status will be returned as usual. (But note that if there were any conditional `INSTEAD` rules, the negation of their qualifications will have been added to the original query. This might reduce the number of rows it processes, and if so the reported status will be affected.)
- If there is any unconditional `INSTEAD` rule for the query, then the original query will not be executed at all. In this case, the server will return the command status for the last query that was inserted by an `INSTEAD` rule (conditional or unconditional) and is of the same command type (`INSERT`, `UPDATE`, or `DELETE`) as the original query. If no query meeting those requirements is added by any rule, then the returned command status shows the original query type and zeroes for the row-count and OID fields.

The programmer can ensure that any desired `INSTEAD` rule is the one that sets the command status in the second case, by giving it the alphabetically last rule name among the active rules, so that it gets applied last.

42.7. Rules Versus Triggers

Many things that can be done using triggers can also be implemented using the Postgres Pro rule system. One of the things that cannot be implemented by rules are some kinds of constraints, especially foreign keys. It is possible to place a qualified rule that rewrites a command to `NOTHING` if the value of a column does not appear in another table. But then the data is silently thrown away and that's not a good idea. If checks for valid values are required, and in the case of an invalid value an error message should be generated, it must be done by a trigger.

In this chapter, we focused on using rules to update views. All of the update rule examples in this chapter can also be implemented using `INSTEAD OF` triggers on the views. Writing such triggers is often easier than writing rules, particularly if complex logic is required to perform the update.

For the things that can be implemented by both, which is best depends on the usage of the database. A trigger is fired once for each affected row. A rule modifies the query or generates an additional query. So if many rows are affected in one statement, a rule issuing one extra command is likely to be faster than a trigger that is called for every single row and must re-determine what to do many times. However, the trigger approach is conceptually far simpler than the rule approach, and is easier for novices to get right.

Here we show an example of how the choice of rules versus triggers plays out in one situation. There are two tables:

```
CREATE TABLE computer (
    hostname      text,      -- indexed
    manufacturer  text      -- indexed
);

CREATE TABLE software (
    software      text,      -- indexed
    hostname      text      -- indexed
);
```

Both tables have many thousands of rows and the indexes on `hostname` are unique. The rule or trigger should implement a constraint that deletes rows from `software` that reference a deleted computer. The trigger would use this command:

```
DELETE FROM software WHERE hostname = $1;
```

Since the trigger is called for each individual row deleted from `computer`, it can prepare and save the plan for this command and pass the `hostname` value in the parameter. The rule would be written as:

```
CREATE RULE computer_del AS ON DELETE TO computer
DO DELETE FROM software WHERE hostname = OLD.hostname;
```

Now we look at different types of deletes. In the case of a:

```
DELETE FROM computer WHERE hostname = 'mypc.local.net';
```

the table `computer` is scanned by index (fast), and the command issued by the trigger would also use an index scan (also fast). The extra command from the rule would be:

```
DELETE FROM software WHERE computer.hostname = 'mypc.local.net'
AND software.hostname = computer.hostname;
```

Since there are appropriate indexes set up, the planner will create a plan of

```
Nestloop
->  Index Scan using comp_hostidx on computer
->  Index Scan using soft_hostidx on software
```

So there would be not that much difference in speed between the trigger and the rule implementation.

With the next delete we want to get rid of all the 2000 computers where the `hostname` starts with `old`. There are two possible commands to do that. One is:

```
DELETE FROM computer WHERE hostname >= 'old'
                        AND hostname < 'ole'
```

The command added by the rule will be:

```
DELETE FROM software WHERE computer.hostname >= 'old' AND computer.hostname < 'ole'
                        AND software.hostname = computer.hostname;
```

with the plan

```
Hash Join
-> Seq Scan on software
-> Hash
    -> Index Scan using comp_hostidx on computer
```

The other possible command is:

```
DELETE FROM computer WHERE hostname ~ '^old';
```

which results in the following executing plan for the command added by the rule:

```
Nestloop
-> Index Scan using comp_hostidx on computer
-> Index Scan using soft_hostidx on software
```

This shows, that the planner does not realize that the qualification for `hostname` in `computer` could also be used for an index scan on `software` when there are multiple qualification expressions combined with `AND`, which is what it does in the regular-expression version of the command. The trigger will get invoked once for each of the 2000 old computers that have to be deleted, and that will result in one index scan over `computer` and 2000 index scans over `software`. The rule implementation will do it with two commands that use indexes. And it depends on the overall size of the table `software` whether the rule will still be faster in the sequential scan situation. 2000 command executions from the trigger over the SPI manager take some time, even if all the index blocks will soon be in the cache.

The last command we look at is:

```
DELETE FROM computer WHERE manufacturer = 'bim';
```

Again this could result in many rows to be deleted from `computer`. So the trigger will again run many commands through the executor. The command generated by the rule will be:

```
DELETE FROM software WHERE computer.manufacturer = 'bim'
                        AND software.hostname = computer.hostname;
```

The plan for that command will again be the nested loop over two index scans, only using a different index on `computer`:

```
Nestloop
-> Index Scan using comp_manufidx on computer
-> Index Scan using soft_hostidx on software
```

In any of these cases, the extra commands from the rule system will be more or less independent from the number of affected rows in a command.

The summary is, rules will only be significantly slower than triggers if their actions result in large and badly qualified joins, a situation where the planner fails.

Chapter 43. Procedural Languages

Postgres Pro allows user-defined functions to be written in other languages besides SQL and C. These other languages are generically called *procedural languages* (PLs). For a function written in a procedural language, the database server has no built-in knowledge about how to interpret the function's source text. Instead, the task is passed to a special handler that knows the details of the language. The handler could either do all the work of parsing, syntax analysis, execution, etc. itself, or it could serve as “glue” between Postgres Pro and an existing implementation of a programming language. The handler itself is a C language function compiled into a shared object and loaded on demand, just like any other C function.

There are currently four procedural languages available in the standard Postgres Pro distribution: PL/pgSQL ([Chapter 44](#)), PL/Tcl ([Chapter 45](#)), PL/Perl ([Chapter 46](#)), and PL/Python ([Chapter 47](#)). There are additional procedural languages available that are not included in the core distribution. [Appendix I](#) has information about finding them. In addition other languages can be defined by users; the basics of developing a new procedural language are covered in [Chapter 57](#).

43.1. Installing Procedural Languages

A procedural language must be “installed” into each database where it is to be used. But procedural languages installed in the database `template1` are automatically available in all subsequently created databases, since their entries in `template1` will be copied by `CREATE DATABASE`. So the database administrator can decide which languages are available in which databases and can make some languages available by default if desired.

For the languages supplied with the standard distribution, it is only necessary to execute `CREATE EXTENSION language_name` to install the language into the current database. The manual procedure described below is only recommended for installing languages that have not been packaged as extensions.

Manual Procedural Language Installation

A procedural language is installed in a database in five steps, which must be carried out by a database superuser. In most cases the required SQL commands should be packaged as the installation script of an “extension”, so that `CREATE EXTENSION` can be used to execute them.

1. The shared object for the language handler must be compiled and installed into an appropriate library directory. This works in the same way as building and installing modules with regular user-defined C functions does; see [Section 39.10.5](#). Often, the language handler will depend on an external library that provides the actual programming language engine; if so, that must be installed as well.
2. The handler must be declared with the command

```
CREATE FUNCTION handler_function_name()  
    RETURNS language_handler  
    AS 'path-to-shared-object'  
    LANGUAGE C;
```

The special return type of `language_handler` tells the database system that this function does not return one of the defined SQL data types and is not directly usable in SQL statements.

3. (Optional) Optionally, the language handler can provide an “inline” handler function that executes anonymous code blocks (`DO` commands) written in this language. If an inline handler function is provided by the language, declare it with a command like

```
CREATE FUNCTION inline_function_name(internal)  
    RETURNS void  
    AS 'path-to-shared-object'  
    LANGUAGE C;
```

4. (Optional) Optionally, the language handler can provide a “validator” function that checks a function definition for correctness without actually executing it. The validator function is called by `CREATE FUNCTION` if it exists. If a validator function is provided by the language, declare it with a command like

```
CREATE FUNCTION validator_function_name(oid)
  RETURNS void
  AS 'path-to-shared-object'
  LANGUAGE C STRICT;
```

5. Finally, the PL must be declared with the command

```
CREATE [TRUSTED] LANGUAGE language_name
  HANDLER handler_function_name
  [INLINE inline_function_name]
  [VALIDATOR validator_function_name] ;
```

The optional key word `TRUSTED` specifies that the language does not grant access to data that the user would not otherwise have. Trusted languages are designed for ordinary database users (those without superuser privilege) and allows them to safely create functions and procedures. Since PL functions are executed inside the database server, the `TRUSTED` flag should only be given for languages that do not allow access to database server internals or the file system. The languages PL/pgSQL, PL/Tcl, and PL/Perl are considered trusted; the languages PL/TclU, PL/PerlU, and PL/PythonU are designed to provide unlimited functionality and should *not* be marked trusted.

[Example 43.1](#) shows how the manual installation procedure would work with the language PL/Perl.

Example 43.1. Manual Installation of PL/Perl

The following command tells the database server where to find the shared object for the PL/Perl language's call handler function:

```
CREATE FUNCTION plperl_call_handler() RETURNS language_handler AS
  '$libdir/plperl' LANGUAGE C;
```

PL/Perl has an inline handler function and a validator function, so we declare those too:

```
CREATE FUNCTION plperl_inline_handler(internal) RETURNS void AS
  '$libdir/plperl' LANGUAGE C STRICT;
```

```
CREATE FUNCTION plperl_validator(oid) RETURNS void AS
  '$libdir/plperl' LANGUAGE C STRICT;
```

The command:

```
CREATE TRUSTED LANGUAGE plperl
  HANDLER plperl_call_handler
  INLINE plperl_inline_handler
  VALIDATOR plperl_validator;
```

then defines that the previously declared functions should be invoked for functions and procedures where the language attribute is `plperl`.

In a default Postgres Pro installation, the handler for the PL/pgSQL language is built and installed into the “library” directory; furthermore, the PL/pgSQL language itself is installed in all databases. If Tcl support is configured in, the handlers for PL/Tcl and PL/TclU are built and installed in the library directory, but the language itself is not installed in any database by default. Likewise, the PL/Perl and PL/PerlU handlers are built and installed if Perl support is configured, and the PL/PythonU handler is installed if Python support is configured, but these languages are not installed by default.

Chapter 44. PL/pgSQL — SQL Procedural Language

44.1. Overview

PL/pgSQL is a loadable procedural language for the Postgres Pro database system. The design goals of PL/pgSQL were to create a loadable procedural language that

- can be used to create functions, procedures, and triggers,
- adds control structures to the SQL language,
- can perform complex computations,
- inherits all user-defined types, functions, procedures, and operators,
- can be defined to be trusted by the server,
- is easy to use.

Functions created with PL/pgSQL can be used anywhere that built-in functions could be used. For example, it is possible to create complex conditional computation functions and later use them to define operators or use them in index expressions.

In PostgreSQL 9.0 and later, PL/pgSQL is installed by default. However it is still a loadable module, so especially security-conscious administrators could choose to remove it.

44.1.1. Advantages of Using PL/pgSQL

SQL is the language Postgres Pro and most other relational databases use as query language. It's portable and easy to learn. But every SQL statement must be executed individually by the database server.

That means that your client application must send each query to the database server, wait for it to be processed, receive and process the results, do some computation, then send further queries to the server. All this incurs interprocess communication and will also incur network overhead if your client is on a different machine than the database server.

With PL/pgSQL you can group a block of computation and a series of queries *inside* the database server, thus having the power of a procedural language and the ease of use of SQL, but with considerable savings of client/server communication overhead.

- Extra round trips between client and server are eliminated
- Intermediate results that the client does not need do not have to be marshaled or transferred between server and client
- Multiple rounds of query parsing can be avoided

This can result in a considerable performance increase as compared to an application that does not use stored functions.

Also, with PL/pgSQL you can use all the data types, operators and functions of SQL.

44.1.2. Supported Argument and Result Data Types

Functions written in PL/pgSQL can accept as arguments any scalar or array data type supported by the server, and they can return a result of any of these types. They can also accept or return any composite type (row type) specified by name. It is also possible to declare a PL/pgSQL function as accepting `record`,

which means that any composite type will do as input, or as returning `record`, which means that the result is a row type whose columns are determined by specification in the calling query, as discussed in [Section 7.2.1.4](#).

PL/pgSQL functions can be declared to accept a variable number of arguments by using the `VARIADIC` marker. This works exactly the same way as for SQL functions, as discussed in [Section 39.5.6](#).

PL/pgSQL functions can also be declared to accept and return the polymorphic types described in [Section 39.2.5](#), thus allowing the actual data types handled by the function to vary from call to call. Examples appear in [Section 44.3.1](#).

PL/pgSQL functions can also be declared to return a “set” (or table) of any data type that can be returned as a single instance. Such a function generates its output by executing `RETURN NEXT` for each desired element of the result set, or by using `RETURN QUERY` to output the result of evaluating a query.

Finally, a PL/pgSQL function can be declared to return `void` if it has no useful return value. (Alternatively, it could be written as a procedure in that case.)

PL/pgSQL functions can also be declared with output parameters in place of an explicit specification of the return type. This does not add any fundamental capability to the language, but it is often convenient, especially for returning multiple values. The `RETURNS TABLE` notation can also be used in place of `RETURNS SETOF`.

Specific examples appear in [Section 44.3.1](#) and [Section 44.6.1](#).

44.2. Structure of PL/pgSQL

Functions written in PL/pgSQL are defined to the server by executing `CREATE FUNCTION` commands. Such a command would normally look like, say,

```
CREATE FUNCTION somefunc(integer, text) RETURNS integer
AS 'function body text'
LANGUAGE plpgsql;
```

The function body is simply a string literal so far as `CREATE FUNCTION` is concerned. It is often helpful to use dollar quoting (see [Section 4.1.2.4](#)) to write the function body, rather than the normal single quote syntax. Without dollar quoting, any single quotes or backslashes in the function body must be escaped by doubling them. Almost all the examples in this chapter use dollar-quoted literals for their function bodies.

PL/pgSQL is a block-structured language. The complete text of a function body must be a *block*. A block is defined as:

```
[ <<label>> ]
[ DECLARE
    declarations ]
BEGIN
    statements
END [ label ];
```

Each declaration and each statement within a block is terminated by a semicolon. A block that appears within another block must have a semicolon after `END`, as shown above; however the final `END` that concludes a function body does not require a semicolon.

Tip

A common mistake is to write a semicolon immediately after `BEGIN`. This is incorrect and will result in a syntax error.

A *label* is only needed if you want to identify the block for use in an `EXIT` statement, or to qualify the names of the variables declared in the block. If a label is given after `END`, it must match the label at the block's beginning.

All key words are case-insensitive. Identifiers are implicitly converted to lower case unless double-quoted, just as they are in ordinary SQL commands.

Comments work the same way in PL/pgSQL code as in ordinary SQL. A double dash (`--`) starts a comment that extends to the end of the line. A `/*` starts a block comment that extends to the matching occurrence of `*/`. Block comments nest.

Any statement in the statement section of a block can be a *subblock*. Subblocks can be used for logical grouping or to localize variables to a small group of statements. Variables declared in a subblock mask any similarly-named variables of outer blocks for the duration of the subblock; but you can access the outer variables anyway if you qualify their names with their block's label. For example:

```
CREATE FUNCTION somefunc() RETURNS integer AS $$
<< outerblock >>
DECLARE
    quantity integer := 30;
BEGIN
    RAISE NOTICE 'Quantity here is %', quantity; -- Prints 30
    quantity := 50;
    --
    -- Create a subblock
    --
    DECLARE
        quantity integer := 80;
    BEGIN
        RAISE NOTICE 'Quantity here is %', quantity; -- Prints 80
        RAISE NOTICE 'Outer quantity here is %', outerblock.quantity; -- Prints 50
    END;

    RAISE NOTICE 'Quantity here is %', quantity; -- Prints 50

    RETURN quantity;
END;
$$ LANGUAGE plpgsql;
```

Note

There is actually a hidden “outer block” surrounding the body of any PL/pgSQL function. This block provides the declarations of the function's parameters (if any), as well as some special variables such as `FOUND` (see [Section 44.5.5](#)). The outer block is labeled with the function's name, meaning that parameters and special variables can be qualified with the function's name.

It is important not to confuse the use of `BEGIN/END` for grouping statements in PL/pgSQL with the similarly-named SQL commands for transaction control. PL/pgSQL's `BEGIN/END` are only for grouping; they do not start or end a transaction. See [Section 44.8](#) for information on managing transactions in PL/pgSQL. Also, a block containing an `EXCEPTION` clause effectively forms a subtransaction that can be rolled back without affecting the outer transaction. For more about that see [Section 44.6.8](#).

44.3. Declarations

All variables used in a block must be declared in the declarations section of the block. (The only exceptions are that the loop variable of a `FOR` loop iterating over a range of integer values is automatically declared as an integer variable, and likewise the loop variable of a `FOR` loop iterating over a cursor's result is automatically declared as a record variable.)

PL/pgSQL variables can have any SQL data type, such as `integer`, `varchar`, and `char`.

Here are some examples of variable declarations:

```
user_id integer;
quantity numeric(5);
url varchar;
myrow tablename%ROWTYPE;
myfield tablename.columnname%TYPE;
arow RECORD;
```

The general syntax of a variable declaration is:

```
name [ CONSTANT ] type [ COLLATE collation_name ] [ NOT NULL ] [ { DEFAULT | := |
= } expression ];
```

The `DEFAULT` clause, if given, specifies the initial value assigned to the variable when the block is entered. If the `DEFAULT` clause is not given then the variable is initialized to the SQL null value. The `CONSTANT` option prevents the variable from being assigned to after initialization, so that its value will remain constant for the duration of the block. The `COLLATE` option specifies a collation to use for the variable (see [Section 44.3.6](#)). If `NOT NULL` is specified, an assignment of a null value results in a run-time error. All variables declared as `NOT NULL` must have a nonnull default value specified. Equal (=) can be used instead of PL/SQL-compliant :=.

A variable's default value is evaluated and assigned to the variable each time the block is entered (not just once per function call). So, for example, assigning `now()` to a variable of type `timestamp` causes the variable to have the time of the current function call, not the time when the function was precompiled.

Examples:

```
quantity integer DEFAULT 32;
url varchar := 'http://mysite.com';
transaction_time CONSTANT timestamp with time zone := now();
```

Once declared, a variable's value can be used in later initialization expressions in the same block, for example:

```
DECLARE
    x integer := 1;
    y integer := x + 1;
```

44.3.1. Declaring Function Parameters

Parameters passed to functions are named with the identifiers `$1`, `$2`, etc. Optionally, aliases can be declared for `$n` parameter names for increased readability. Either the alias or the numeric identifier can then be used to refer to the parameter value.

There are two ways to create an alias. The preferred way is to give a name to the parameter in the `CREATE FUNCTION` command, for example:

```
CREATE FUNCTION sales_tax(subtotal real) RETURNS real AS $$
BEGIN
    RETURN subtotal * 0.06;
END;
$$ LANGUAGE plpgsql;
```

The other way is to explicitly declare an alias, using the declaration syntax

```
name ALIAS FOR $n;
```

The same example in this style looks like:

```
CREATE FUNCTION sales_tax(real) RETURNS real AS $$
```

```
DECLARE
    subtotal ALIAS FOR $1;
BEGIN
    RETURN subtotal * 0.06;
END;
$$ LANGUAGE plpgsql;
```

Note

These two examples are not perfectly equivalent. In the first case, `subtotal` could be referenced as `sales_tax.subtotal`, but in the second case it could not. (Had we attached a label to the inner block, `subtotal` could be qualified with that label, instead.)

Some more examples:

```
CREATE FUNCTION instr(vchar, integer) RETURNS integer AS $$
DECLARE
    v_string ALIAS FOR $1;
    index ALIAS FOR $2;
BEGIN
    -- some computations using v_string and index here
END;
$$ LANGUAGE plpgsql;
```

```
CREATE FUNCTION concat_selected_fields(in_t sometablename) RETURNS text AS $$
BEGIN
    RETURN in_t.f1 || in_t.f3 || in_t.f5 || in_t.f7;
END;
$$ LANGUAGE plpgsql;
```

When a PL/pgSQL function is declared with output parameters, the output parameters are given n names and optional aliases in just the same way as the normal input parameters. An output parameter is effectively a variable that starts out NULL; it should be assigned to during the execution of the function. The final value of the parameter is what is returned. For instance, the sales-tax example could also be done this way:

```
CREATE FUNCTION sales_tax(subtotal real, OUT tax real) AS $$
BEGIN
    tax := subtotal * 0.06;
END;
$$ LANGUAGE plpgsql;
```

Notice that we omitted `RETURNS real` — we could have included it, but it would be redundant.

To call a function with `OUT` parameters, omit the output parameter(s) in the function call:

```
SELECT sales_tax(100.00);
```

Output parameters are most useful when returning multiple values. A trivial example is:

```
CREATE FUNCTION sum_n_product(x int, y int, OUT sum int, OUT prod int) AS $$
BEGIN
    sum := x + y;
    prod := x * y;
END;
$$ LANGUAGE plpgsql;
```

```
SELECT * FROM sum_n_product(2, 4);
sum | prod
-----+-----
6 |      8
```

As discussed in [Section 39.5.4](#), this effectively creates an anonymous record type for the function's results. If a `RETURNS` clause is given, it must say `RETURNS record`.

This also works with procedures, for example:

```
CREATE PROCEDURE sum_n_product(x int, y int, OUT sum int, OUT prod int) AS $$
BEGIN
    sum := x + y;
    prod := x * y;
END;
$$ LANGUAGE plpgsql;
```

In a call to a procedure, all the parameters must be specified. For output parameters, `NULL` may be specified when calling the procedure from plain SQL:

```
CALL sum_n_product(2, 4, NULL, NULL);
sum | prod
-----+-----
6 |      8
```

However, when calling a procedure from PL/pgSQL, you should instead write a variable for any output parameter; the variable will receive the result of the call. See [Section 44.6.3](#) for details.

Another way to declare a PL/pgSQL function is with `RETURNS TABLE`, for example:

```
CREATE FUNCTION extended_sales(p_itemno int)
RETURNS TABLE(quantity int, total numeric) AS $$
BEGIN
    RETURN QUERY SELECT s.quantity, s.quantity * s.price FROM sales AS s
                  WHERE s.itemno = p_itemno;
END;
$$ LANGUAGE plpgsql;
```

This is exactly equivalent to declaring one or more `OUT` parameters and specifying `RETURNS SETOF some-type`.

When the return type of a PL/pgSQL function is declared as a polymorphic type (see [Section 39.2.5](#)), a special parameter `$0` is created. Its data type is the actual return type of the function, as deduced from the actual input types. This allows the function to access its actual return type as shown in [Section 44.3.3](#). `$0` is initialized to null and can be modified by the function, so it can be used to hold the return value if desired, though that is not required. `$0` can also be given an alias. For example, this function works on any data type that has a `+` operator:

```
CREATE FUNCTION add_three_values(v1 anyelement, v2 anyelement, v3 anyelement)
RETURNS anyelement AS $$
DECLARE
    result ALIAS FOR $0;
BEGIN
    result := v1 + v2 + v3;
    RETURN result;
END;
$$ LANGUAGE plpgsql;
```

The same effect can be obtained by declaring one or more output parameters as polymorphic types. In this case the special `$0` parameter is not used; the output parameters themselves serve the same purpose. For example:

```
CREATE FUNCTION add_three_values(v1 anyelement, v2 anyelement, v3 anyelement,
                                OUT sum anyelement)
AS $$
BEGIN
    sum := v1 + v2 + v3;
END;
$$ LANGUAGE plpgsql;
```

In practice it might be more useful to declare a polymorphic function using the `anycompatible` family of types, so that automatic promotion of the input arguments to a common type will occur. For example:

```
CREATE FUNCTION add_three_values(v1 anycompatible, v2 anycompatible, v3 anycompatible)
RETURNS anycompatible AS $$
BEGIN
    RETURN v1 + v2 + v3;
END;
$$ LANGUAGE plpgsql;
```

With this example, a call such as

```
SELECT add_three_values(1, 2, 4.7);
```

will work, automatically promoting the integer inputs to numeric. The function using `anyelement` would require you to cast the three inputs to the same type manually.

44.3.2. ALIAS

```
newname ALIAS FOR oldname;
```

The `ALIAS` syntax is more general than is suggested in the previous section: you can declare an alias for any variable, not just function parameters. The main practical use for this is to assign a different name for variables with predetermined names, such as `NEW` or `OLD` within a trigger function.

Examples:

```
DECLARE
    prior ALIAS FOR old;
    updated ALIAS FOR new;
```

Since `ALIAS` creates two different ways to name the same object, unrestricted use can be confusing. It's best to use it only for the purpose of overriding predetermined names.

44.3.3. Copying Types

```
name table.column%TYPE
name variable%TYPE
```

`%TYPE` provides the data type of a table column or a previously-declared PL/pgSQL variable. You can use this to declare variables that will hold database values. For example, let's say you have a column named `user_id` in your `users` table. To declare a variable with the same data type as `users.user_id` you write:

```
user_id users.user_id%TYPE;
```

It is also possible to write array decoration after `%TYPE`, thereby creating a variable that holds an array of the referenced type:

```
user_ids users.user_id%TYPE[];
user_ids users.user_id%TYPE ARRAY[4]; -- equivalent to the above
```

Just as when declaring table columns that are arrays, it doesn't matter whether you write multiple bracket pairs or specific array dimensions: Postgres Pro treats all arrays of a given element type as the same type, regardless of dimensionality. (See [Section 8.15.1](#).)

By using `%TYPE` you don't need to know the data type of the structure you are referencing, and most importantly, if the data type of the referenced item changes in the future (for instance: you change the type of `user_id` from `integer` to `real`), you might not need to change your function definition.

`%TYPE` is particularly valuable in polymorphic functions, since the data types needed for internal variables can change from one call to the next. Appropriate variables can be created by applying `%TYPE` to the function's arguments or result placeholders.

44.3.4. Row Types

```
name table_name%ROWTYPE;
name composite_type_name;
```

A variable of a composite type is called a *row variable* (or *row-type variable*). Such a variable can hold a whole row of a `SELECT` or `FOR` query result, so long as that query's column set matches the declared type of the variable. The individual fields of the row value are accessed using the usual dot notation, for example `rowvar.field`.

A row variable can be declared to have the same type as the rows of an existing table or view, by using the `table_name%ROWTYPE` notation; or it can be declared by giving a composite type's name. (Since every table has an associated composite type of the same name, it actually does not matter in Postgres Pro whether you write `%ROWTYPE` or not. But the form with `%ROWTYPE` is more portable.)

As with `%TYPE`, `%ROWTYPE` can be followed by array decoration to declare a variable that holds an array of the referenced composite type.

Parameters to a function can be composite types (complete table rows). In that case, the corresponding identifier `$n` will be a row variable, and fields can be selected from it, for example `$1.user_id`.

Here is an example of using composite types. `table1` and `table2` are existing tables having at least the mentioned fields:

```
CREATE FUNCTION merge_fields(t_row table1) RETURNS text AS $$
DECLARE
    t2_row table2%ROWTYPE;
BEGIN
    SELECT * INTO t2_row FROM table2 WHERE ... ;
    RETURN t_row.f1 || t2_row.f3 || t_row.f5 || t2_row.f7;
END;
$$ LANGUAGE plpgsql;

SELECT merge_fields(t.*) FROM table1 t WHERE ... ;
```

44.3.5. Record Types

```
name RECORD;
```

Record variables are similar to row-type variables, but they have no predefined structure. They take on the actual row structure of the row they are assigned during a `SELECT` or `FOR` command. The substructure of a record variable can change each time it is assigned to. A consequence of this is that until a record variable is first assigned to, it has no substructure, and any attempt to access a field in it will draw a run-time error.

Note that `RECORD` is not a true data type, only a placeholder. One should also realize that when a PL/pgSQL function is declared to return type `record`, this is not quite the same concept as a record variable, even though such a function might use a record variable to hold its result. In both cases the actual row structure is unknown when the function is written, but for a function returning `record` the actual structure is determined when the calling query is parsed, whereas a record variable can change its row structure on-the-fly.

44.3.6. Collation of PL/pgSQL Variables

When a PL/pgSQL function has one or more parameters of collatable data types, a collation is identified for each function call depending on the collations assigned to the actual arguments, as described in [Section 23.2](#). If a collation is successfully identified (i.e., there are no conflicts of implicit collations among the arguments) then all the collatable parameters are treated as having that collation implicitly. This will affect the behavior of collation-sensitive operations within the function. For example, consider

```
CREATE FUNCTION less_than(a text, b text) RETURNS boolean AS $$
BEGIN
    RETURN a < b;
END;
$$ LANGUAGE plpgsql;

SELECT less_than(text_field_1, text_field_2) FROM table1;
SELECT less_than(text_field_1, text_field_2 COLLATE "C") FROM table1;
```

The first use of `less_than` will use the common collation of `text_field_1` and `text_field_2` for the comparison, while the second use will use C collation.

Furthermore, the identified collation is also assumed as the collation of any local variables that are of collatable types. Thus this function would not work any differently if it were written as

```
CREATE FUNCTION less_than(a text, b text) RETURNS boolean AS $$
DECLARE
    local_a text := a;
    local_b text := b;
BEGIN
    RETURN local_a < local_b;
END;
$$ LANGUAGE plpgsql;
```

If there are no parameters of collatable data types, or no common collation can be identified for them, then parameters and local variables use the default collation of their data type (which is usually the database's default collation, but could be different for variables of domain types).

A local variable of a collatable data type can have a different collation associated with it by including the `COLLATE` option in its declaration, for example

```
DECLARE
    local_a text COLLATE "en_US";
```

This option overrides the collation that would otherwise be given to the variable according to the rules above.

Also, of course explicit `COLLATE` clauses can be written inside a function if it is desired to force a particular collation to be used in a particular operation. For example,

```
CREATE FUNCTION less_than_c(a text, b text) RETURNS boolean AS $$
BEGIN
    RETURN a < b COLLATE "C";
END;
$$ LANGUAGE plpgsql;
```

This overrides the collations associated with the table columns, parameters, or local variables used in the expression, just as would happen in a plain SQL command.

44.4. Expressions

All expressions used in PL/pgSQL statements are processed using the server's main SQL executor. For example, when you write a PL/pgSQL statement like

```
IF expression THEN ...
```

PL/pgSQL will evaluate the expression by feeding a query like

```
SELECT expression
```

to the main SQL engine. While forming the `SELECT` command, any occurrences of PL/pgSQL variable names are replaced by query parameters, as discussed in detail in [Section 44.12.1](#). This allows the query plan for the `SELECT` to be prepared just once and then reused for subsequent evaluations with different values of the variables. Thus, what really happens on first use of an expression is essentially a `PREPARE` command. For example, if we have declared two integer variables `x` and `y`, and we write

```
IF x < y THEN ...
```

what happens behind the scenes is equivalent to

```
PREPARE statement_name(integer, integer) AS SELECT $1 < $2;
```

and then this prepared statement is `EXECUTED` for each execution of the `IF` statement, with the current values of the PL/pgSQL variables supplied as parameter values. Normally these details are not important to a PL/pgSQL user, but they are useful to know when trying to diagnose a problem. More information appears in [Section 44.12.2](#).

Since an *expression* is converted to a `SELECT` command, it can contain the same clauses that an ordinary `SELECT` would, except that it cannot include a top-level `UNION`, `INTERSECT`, or `EXCEPT` clause. Thus for example one could test whether a table is non-empty with

```
IF count(*) > 0 FROM my_table THEN ...
```

since the *expression* between `IF` and `THEN` is parsed as though it were `SELECT count(*) > 0 FROM my_table`. The `SELECT` must produce a single column, and not more than one row. (If it produces no rows, the result is taken as `NULL`.)

44.5. Basic Statements

In this section and the following ones, we describe all the statement types that are explicitly understood by PL/pgSQL. Anything not recognized as one of these statement types is presumed to be an SQL command and is sent to the main database engine to execute, as described in [Section 44.5.2](#).

44.5.1. Assignment

An assignment of a value to a PL/pgSQL variable is written as:

```
variable { := | = } expression;
```

As explained previously, the expression in such a statement is evaluated by means of an SQL `SELECT` command sent to the main database engine. The expression must yield a single value (possibly a row value, if the variable is a row or record variable). The target variable can be a simple variable (optionally qualified with a block name), a field of a row or record target, or an element or slice of an array target. Equal (`=`) can be used instead of PL/SQL-compliant `:=`.

If the expression's result data type doesn't match the variable's data type, the value will be coerced as though by an assignment cast (see [Section 10.4](#)). If no assignment cast is known for the pair of data types involved, the PL/pgSQL interpreter will attempt to convert the result value textually, that is by applying the result type's output function followed by the variable type's input function. Note that this could result in run-time errors generated by the input function, if the string form of the result value is not acceptable to the input function.

Examples:

```
tax := subtotal * 0.06;
my_record.user_id := 20;
my_array[j] := 20;
my_array[1:3] := array[1,2,3];
```

```
complex_array[n].realpart = 12.3;
```

44.5.2. Executing SQL Commands

In general, any SQL command that does not return rows can be executed within a PL/pgSQL function just by writing the command. For example, you could create and fill a table by writing

```
CREATE TABLE mytable (id int primary key, data text);
INSERT INTO mytable VALUES (1, 'one'), (2, 'two');
```

If the command does return rows (for example `SELECT`, or `INSERT/UPDATE/DELETE/MERGE` with `RETURNING`), there are two ways to proceed. When the command will return at most one row, or you only care about the first row of output, write the command as usual but add an `INTO` clause to capture the output, as described in [Section 44.5.3](#). To process all of the output rows, write the command as the data source for a `FOR` loop, as described in [Section 44.6.6](#).

Usually it is not sufficient just to execute statically-defined SQL commands. Typically you'll want a command to use varying data values, or even to vary in more fundamental ways such as by using different table names at different times. Again, there are two ways to proceed depending on the situation.

PL/pgSQL variable values can be automatically inserted into optimizable SQL commands, which are `SELECT`, `INSERT`, `UPDATE`, `DELETE`, `MERGE`, and certain utility commands that incorporate one of these, such as `EXPLAIN` and `CREATE TABLE ... AS SELECT`. In these commands, any PL/pgSQL variable name appearing in the command text is replaced by a query parameter, and then the current value of the variable is provided as the parameter value at run time. This is exactly like the processing described earlier for expressions; for details see [Section 44.12.1](#).

When executing an optimizable SQL command in this way, PL/pgSQL may cache and re-use the execution plan for the command, as discussed in [Section 44.12.2](#).

Non-optimizable SQL commands (also called utility commands) are not capable of accepting query parameters. So automatic substitution of PL/pgSQL variables does not work in such commands. To include non-constant text in a utility command executed from PL/pgSQL, you must build the utility command as a string and then `EXECUTE` it, as discussed in [Section 44.5.4](#).

`EXECUTE` must also be used if you want to modify the command in some other way than supplying a data value, for example by changing a table name.

Sometimes it is useful to evaluate an expression or `SELECT` query but discard the result, for example when calling a function that has side-effects but no useful result value. To do this in PL/pgSQL, use the `PERFORM` statement:

```
PERFORM query;
```

This executes *query* and discards the result. Write the *query* the same way you would write an SQL `SELECT` command, but replace the initial keyword `SELECT` with `PERFORM`. For `WITH` queries, use `PERFORM` and then place the query in parentheses. (In this case, the query can only return one row.) PL/pgSQL variables will be substituted into the query just as described above, and the plan is cached in the same way. Also, the special variable `FOUND` is set to true if the query produced at least one row, or false if it produced no rows (see [Section 44.5.5](#)).

Note

One might expect that writing `SELECT` directly would accomplish this result, but at present the only accepted way to do it is `PERFORM`. An SQL command that can return rows, such as `SELECT`, will be rejected as an error unless it has an `INTO` clause as discussed in the next section.

An example:


```
PERFORM create_mv('cs_session_page_requests_mv', my_query);
```

44.5.3. Executing a Command with a Single-Row Result

The result of an SQL command yielding a single row (possibly of multiple columns) can be assigned to a record variable, row-type variable, or list of scalar variables. This is done by writing the base SQL command and adding an `INTO` clause. For example,

```
SELECT select_expressions INTO [STRICT] target FROM ...;
INSERT ... RETURNING expressions INTO [STRICT] target;
UPDATE ... RETURNING expressions INTO [STRICT] target;
DELETE ... RETURNING expressions INTO [STRICT] target;
MERGE ... RETURNING expressions INTO [STRICT] target;
```

where *target* can be a record variable, a row variable, or a comma-separated list of simple variables and record/row fields. PL/pgSQL variables will be substituted into the rest of the command (that is, everything but the `INTO` clause) just as described above, and the plan is cached in the same way. This works for `SELECT`, `INSERT/UPDATE/DELETE/MERGE` with `RETURNING`, and certain utility commands that return row sets, such as `EXPLAIN`. Except for the `INTO` clause, the SQL command is the same as it would be written outside PL/pgSQL.

Tip

Note that this interpretation of `SELECT` with `INTO` is quite different from Postgres Pro's regular `SELECT INTO` command, wherein the `INTO` target is a newly created table. If you want to create a table from a `SELECT` result inside a PL/pgSQL function, use the syntax `CREATE TABLE ... AS SELECT`.

If a row variable or a variable list is used as target, the command's result columns must exactly match the structure of the target as to number and data types, or else a run-time error occurs. When a record variable is the target, it automatically configures itself to the row type of the command's result columns.

The `INTO` clause can appear almost anywhere in the SQL command. Customarily it is written either just before or just after the list of *select_expressions* in a `SELECT` command, or at the end of the command for other command types. It is recommended that you follow this convention in case the PL/pgSQL parser becomes stricter in future versions.

If `STRICT` is not specified in the `INTO` clause, then *target* will be set to the first row returned by the command, or to nulls if the command returned no rows. (Note that “the first row” is not well-defined unless you've used `ORDER BY`.) Any result rows after the first row are discarded. You can check the special `FOUND` variable (see [Section 44.5.5](#)) to determine whether a row was returned:

```
SELECT * INTO myrec FROM emp WHERE empname = myname;
IF NOT FOUND THEN
    RAISE EXCEPTION 'employee % not found', myname;
END IF;
```

If the `STRICT` option is specified, the command must return exactly one row or a run-time error will be reported, either `NO_DATA_FOUND` (no rows) or `TOO_MANY_ROWS` (more than one row). You can use an exception block if you wish to catch the error, for example:

```
BEGIN
    SELECT * INTO STRICT myrec FROM emp WHERE empname = myname;
EXCEPTION
    WHEN NO_DATA_FOUND THEN
        RAISE EXCEPTION 'employee % not found', myname;
    WHEN TOO_MANY_ROWS THEN
        RAISE EXCEPTION 'employee % not unique', myname;
END;
```

Successful execution of a command with `STRICT` always sets `FOUND` to true.

For `INSERT/UPDATE/DELETE/MERGE` with `RETURNING`, PL/pgSQL reports an error for more than one returned row, even when `STRICT` is not specified. This is because there is no option such as `ORDER BY` with which to determine which affected row should be returned.

If `print_strict_params` is enabled for the function, then when an error is thrown because the requirements of `STRICT` are not met, the `DETAIL` part of the error message will include information about the parameters passed to the command. You can change the `print_strict_params` setting for all functions by setting `plpgsql.print_strict_params`, though only subsequent function compilations will be affected. You can also enable it on a per-function basis by using a compiler option, for example:

```
CREATE FUNCTION get_userid(username text) RETURNS int
AS $$
#print_strict_params on
DECLARE
userid int;
BEGIN
    SELECT users.userid INTO STRICT userid
        FROM users WHERE users.username = get_userid.username;
    RETURN userid;
END;
$$ LANGUAGE plpgsql;
```

On failure, this function might produce an error message such as

```
ERROR:  query returned no rows
DETAIL:  parameters: username = 'nosuchuser'
CONTEXT:  PL/pgSQL function get_userid(text) line 6 at SQL statement
```

Note

The `STRICT` option matches the behavior of Oracle PL/SQL's `SELECT INTO` and related statements.

44.5.4. Executing Dynamic Commands

Oftentimes you will want to generate dynamic commands inside your PL/pgSQL functions, that is, commands that will involve different tables or different data types each time they are executed. PL/pgSQL's normal attempts to cache plans for commands (as discussed in [Section 44.12.2](#)) will not work in such scenarios. To handle this sort of problem, the `EXECUTE` statement is provided:

```
EXECUTE command-string [ INTO [STRICT] target ] [ USING expression [, ... ] ];
```

where *command-string* is an expression yielding a string (of type `text`) containing the command to be executed. The optional *target* is a record variable, a row variable, or a comma-separated list of simple variables and record/row fields, into which the results of the command will be stored. The optional `USING` expressions supply values to be inserted into the command.

No substitution of PL/pgSQL variables is done on the computed command string. Any required variable values must be inserted in the command string as it is constructed; or you can use parameters as described below.

Also, there is no plan caching for commands executed via `EXECUTE`. Instead, the command is always planned each time the statement is run. Thus the command string can be dynamically created within the function to perform actions on different tables and columns.

The `INTO` clause specifies where the results of an SQL command returning rows should be assigned. If a row variable or variable list is provided, it must exactly match the structure of the command's results; if a record variable is provided, it will configure itself to match the result structure automatically. If multiple

rows are returned, only the first will be assigned to the `INTO` variable(s). If no rows are returned, `NULL` is assigned to the `INTO` variable(s). If no `INTO` clause is specified, the command results are discarded.

If the `STRICT` option is given, an error is reported unless the command produces exactly one row.

The command string can use parameter values, which are referenced in the command as `$1`, `$2`, etc. These symbols refer to values supplied in the `USING` clause. This method is often preferable to inserting data values into the command string as text: it avoids run-time overhead of converting the values to text and back, and it is much less prone to SQL-injection attacks since there is no need for quoting or escaping. An example is:

```
EXECUTE 'SELECT count(*) FROM mytable WHERE inserted_by = $1 AND inserted <= $2'
      INTO c
      USING checked_user, checked_date;
```

Note that parameter symbols can only be used for data values — if you want to use dynamically determined table or column names, you must insert them into the command string textually. For example, if the preceding query needed to be done against a dynamically selected table, you could do this:

```
EXECUTE 'SELECT count(*) FROM '
      || quote_ident(tabname)
      || ' WHERE inserted_by = $1 AND inserted <= $2'
      INTO c
      USING checked_user, checked_date;
```

A cleaner approach is to use `format()`'s `%I` specification to insert table or column names with automatic quoting:

```
EXECUTE format('SELECT count(*) FROM %I '
      'WHERE inserted_by = $1 AND inserted <= $2', tabname)
      INTO c
      USING checked_user, checked_date;
```

(This example relies on the SQL rule that string literals separated by a newline are implicitly concatenated.)

Another restriction on parameter symbols is that they only work in optimizable SQL commands (`SELECT`, `INSERT`, `UPDATE`, `DELETE`, `MERGE`, and certain commands containing one of these). In other statement types (generically called utility statements), you must insert values textually even if they are just data values.

An `EXECUTE` with a simple constant command string and some `USING` parameters, as in the first example above, is functionally equivalent to just writing the command directly in PL/pgSQL and allowing replacement of PL/pgSQL variables to happen automatically. The important difference is that `EXECUTE` will re-plan the command on each execution, generating a plan that is specific to the current parameter values; whereas PL/pgSQL may otherwise create a generic plan and cache it for re-use. In situations where the best plan depends strongly on the parameter values, it can be helpful to use `EXECUTE` to positively ensure that a generic plan is not selected.

`SELECT INTO` is not currently supported within `EXECUTE`; instead, execute a plain `SELECT` command and specify `INTO` as part of the `EXECUTE` itself.

Note

The PL/pgSQL `EXECUTE` statement is not related to the `EXECUTE` SQL statement supported by the Postgres Pro server. The server's `EXECUTE` statement cannot be used directly within PL/pgSQL functions (and is not needed).

Example 44.1. Quoting Values in Dynamic Queries

When working with dynamic commands you will often have to handle escaping of single quotes. The recommended method for quoting fixed text in your function body is dollar quoting. (If you have legacy code that does not use dollar quoting, please refer to the overview in [Section 44.13.1](#), which can save you some effort when translating said code to a more reasonable scheme.)

Dynamic values require careful handling since they might contain quote characters. An example using `format()` (this assumes that you are dollar quoting the function body so quote marks need not be doubled):

```
EXECUTE format('UPDATE tbl SET %I = $1 '
              'WHERE key = $2', colname) USING newvalue, keyvalue;
```

It is also possible to call the quoting functions directly:

```
EXECUTE 'UPDATE tbl SET '
      || quote_ident(colname)
      || ' = '
      || quote_literal(newvalue)
      || ' WHERE key = '
      || quote_literal(keyvalue);
```

This example demonstrates the use of the `quote_ident` and `quote_literal` functions (see [Section 9.4](#)). For safety, expressions containing column or table identifiers should be passed through `quote_ident` before insertion in a dynamic query. Expressions containing values that should be literal strings in the constructed command should be passed through `quote_literal`. These functions take the appropriate steps to return the input text enclosed in double or single quotes respectively, with any embedded special characters properly escaped.

Because `quote_literal` is labeled `STRICT`, it will always return null when called with a null argument. In the above example, if `newvalue` or `keyvalue` were null, the entire dynamic query string would become null, leading to an error from `EXECUTE`. You can avoid this problem by using the `quote_nullable` function, which works the same as `quote_literal` except that when called with a null argument it returns the string `NULL`. For example,

```
EXECUTE 'UPDATE tbl SET '
      || quote_ident(colname)
      || ' = '
      || quote_nullable(newvalue)
      || ' WHERE key = '
      || quote_nullable(keyvalue);
```

If you are dealing with values that might be null, you should usually use `quote_nullable` in place of `quote_literal`.

As always, care must be taken to ensure that null values in a query do not deliver unintended results. For example the `WHERE` clause

```
'WHERE key = ' || quote_nullable(keyvalue)
```

will never succeed if `keyvalue` is null, because the result of using the equality operator `=` with a null operand is always null. If you wish null to work like an ordinary key value, you would need to rewrite the above as

```
'WHERE key IS NOT DISTINCT FROM ' || quote_nullable(keyvalue)
```

(At present, `IS NOT DISTINCT FROM` is handled much less efficiently than `=`, so don't do this unless you must. See [Section 9.2](#) for more information on nulls and `IS DISTINCT`.)

Note that dollar quoting is only useful for quoting fixed text. It would be a very bad idea to try to write this example as:

```
EXECUTE 'UPDATE tbl SET '
      || quote_ident(colname)
```

```
|| ' = $$'
|| newvalue
|| '$$ WHERE key = '
|| quote_literal(keyvalue);
```

because it would break if the contents of `newvalue` happened to contain `$$`. The same objection would apply to any other dollar-quoting delimiter you might pick. So, to safely quote text that is not known in advance, you *must* use `quote_literal`, `quote_nullable`, or `quote_ident`, as appropriate.

Dynamic SQL statements can also be safely constructed using the `format` function (see [Section 9.4.1](#)). For example:

```
EXECUTE format('UPDATE tbl SET %I = %L '
              'WHERE key = %L', colname, newvalue, keyvalue);
```

`%I` is equivalent to `quote_ident`, and `%L` is equivalent to `quote_nullable`. The `format` function can be used in conjunction with the `USING` clause:

```
EXECUTE format('UPDATE tbl SET %I = $1 WHERE key = $2', colname)
      USING newvalue, keyvalue;
```

This form is better because the variables are handled in their native data type format, rather than unconditionally converting them to text and quoting them via `%L`. It is also more efficient.

A much larger example of a dynamic command and `EXECUTE` can be seen in [Example 44.10](#), which builds and executes a `CREATE FUNCTION` command to define a new function.

44.5.5. Obtaining the Result Status

There are several ways to determine the effect of a command. The first method is to use the `GET DIAGNOSTICS` command, which has the form:

```
GET [ CURRENT ] DIAGNOSTICS variable { = | := } item [ , ... ];
```

This command allows retrieval of system status indicators. `CURRENT` is a noise word (but see also `GET STACKED DIAGNOSTICS` in [Section 44.6.8.1](#)). Each *item* is a key word identifying a status value to be assigned to the specified *variable* (which should be of the right data type to receive it). The currently available status items are shown in [Table 44.1](#). Colon-equal (`:=`) can be used instead of the SQL-standard `=` token. An example:

```
GET DIAGNOSTICS integer_var = ROW_COUNT;
```

Table 44.1. Available Diagnostics Items

Name	Type	Description
ROW_COUNT	bigint	the number of rows processed by the most recent SQL command
PG_CONTEXT	text	line(s) of text describing the current call stack (see Section 44.6.9)
PG_ROUTINE_OID	oid	OID of the current function

The second method to determine the effects of a command is to check the special variable named `FOUND`, which is of type `boolean`. `FOUND` starts out false within each PL/pgSQL function call. It is set by each of the following types of statements:

- A `SELECT INTO` statement sets `FOUND` true if a row is assigned, false if no row is returned.
- A `PERFORM` statement sets `FOUND` true if it produces (and discards) one or more rows, false if no row is produced.
- `UPDATE`, `INSERT`, `DELETE`, and `MERGE` statements set `FOUND` true if at least one row is affected, false if no row is affected.

- A `FETCH` statement sets `FOUND` true if it returns a row, false if no row is returned.
- A `MOVE` statement sets `FOUND` true if it successfully repositions the cursor, false otherwise.
- A `FOR` or `FOREACH` statement sets `FOUND` true if it iterates one or more times, else false. `FOUND` is set this way when the loop exits; inside the execution of the loop, `FOUND` is not modified by the loop statement, although it might be changed by the execution of other statements within the loop body.
- `RETURN QUERY` and `RETURN QUERY EXECUTE` statements set `FOUND` true if the query returns at least one row, false if no row is returned.

Other PL/pgSQL statements do not change the state of `FOUND`. Note in particular that `EXECUTE` changes the output of `GET DIAGNOSTICS`, but does not change `FOUND`.

`FOUND` is a local variable within each PL/pgSQL function; any changes to it affect only the current function.

44.5.6. Doing Nothing At All

Sometimes a placeholder statement that does nothing is useful. For example, it can indicate that one arm of an if/then/else chain is deliberately empty. For this purpose, use the `NULL` statement:

```
NULL;
```

For example, the following two fragments of code are equivalent:

```
BEGIN
    y := x / 0;
EXCEPTION
    WHEN division_by_zero THEN
        NULL; -- ignore the error
END;
```

```
BEGIN
    y := x / 0;
EXCEPTION
    WHEN division_by_zero THEN -- ignore the error
END;
```

Which is preferable is a matter of taste.

Note

In Oracle's PL/SQL, empty statement lists are not allowed, and so `NULL` statements are *required* for situations such as this. PL/pgSQL allows you to just write nothing, instead.

44.6. Control Structures

Control structures are probably the most useful (and important) part of PL/pgSQL. With PL/pgSQL's control structures, you can manipulate Postgres Pro data in a very flexible and powerful way.

44.6.1. Returning from a Function

There are two commands available that allow you to return data from a function: `RETURN` and `RETURN NEXT`.

44.6.1.1. RETURN

```
RETURN expression;
```

`RETURN` with an expression terminates the function and returns the value of *expression* to the caller. This form is used for PL/pgSQL functions that do not return a set.

In a function that returns a scalar type, the expression's result will automatically be cast into the function's return type as described for assignments. But to return a composite (row) value, you must write an expression delivering exactly the requested column set. This may require use of explicit casting.

If you declared the function with output parameters, write just `RETURN` with no expression. The current values of the output parameter variables will be returned.

If you declared the function to return `void`, a `RETURN` statement can be used to exit the function early; but do not write an expression following `RETURN`.

The return value of a function cannot be left undefined. If control reaches the end of the top-level block of the function without hitting a `RETURN` statement, a run-time error will occur. This restriction does not apply to functions with output parameters and functions returning `void`, however. In those cases a `RETURN` statement is automatically executed if the top-level block finishes.

Some examples:

```
-- functions returning a scalar type
RETURN 1 + 2;
RETURN scalar_var;

-- functions returning a composite type
RETURN composite_type_var;
RETURN (1, 2, 'three'::text); -- must cast columns to correct types
```

44.6.1.2. RETURN NEXT and RETURN QUERY

```
RETURN NEXT expression;
RETURN QUERY query;
RETURN QUERY EXECUTE command-string [ USING expression [, ... ] ];
```

When a PL/pgSQL function is declared to return `SETOF sometype`, the procedure to follow is slightly different. In that case, the individual items to return are specified by a sequence of `RETURN NEXT` or `RETURN QUERY` commands, and then a final `RETURN` command with no argument is used to indicate that the function has finished executing. `RETURN NEXT` can be used with both scalar and composite data types; with a composite result type, an entire “table” of results will be returned. `RETURN QUERY` appends the results of executing a query to the function's result set. `RETURN NEXT` and `RETURN QUERY` can be freely intermixed in a single set-returning function, in which case their results will be concatenated.

`RETURN NEXT` and `RETURN QUERY` do not actually return from the function — they simply append zero or more rows to the function's result set. Execution then continues with the next statement in the PL/pgSQL function. As successive `RETURN NEXT` or `RETURN QUERY` commands are executed, the result set is built up. A final `RETURN`, which should have no argument, causes control to exit the function (or you can just let control reach the end of the function).

`RETURN QUERY` has a variant `RETURN QUERY EXECUTE`, which specifies the query to be executed dynamically. Parameter expressions can be inserted into the computed query string via `USING`, in just the same way as in the `EXECUTE` command.

If you declared the function with output parameters, write just `RETURN NEXT` with no expression. On each execution, the current values of the output parameter variable(s) will be saved for eventual return as a row of the result. Note that you must declare the function as returning `SETOF record` when there are multiple output parameters, or `SETOF sometype` when there is just one output parameter of type *sometype*, in order to create a set-returning function with output parameters.

Here is an example of a function using `RETURN NEXT`:

```
CREATE TABLE foo (fooid INT, foosubid INT, fooname TEXT);
INSERT INTO foo VALUES (1, 2, 'three');
INSERT INTO foo VALUES (4, 5, 'six');

CREATE OR REPLACE FUNCTION get_all_foo() RETURNS SETOF foo AS
$BODY$
DECLARE
    r foo%rowtype;
BEGIN
    FOR r IN
        SELECT * FROM foo WHERE fooid > 0
    LOOP
        -- can do some processing here
        RETURN NEXT r; -- return current row of SELECT
    END LOOP;
    RETURN;
END;
$BODY$
LANGUAGE plpgsql;

SELECT * FROM get_all_foo();
```

Here is an example of a function using RETURN QUERY:

```
CREATE FUNCTION get_available_flightid(date) RETURNS SETOF integer AS
$BODY$
BEGIN
    RETURN QUERY SELECT flightid
                  FROM flight
                  WHERE flightdate >= $1
                     AND flightdate < ($1 + 1);

    -- Since execution is not finished, we can check whether rows were returned
    -- and raise exception if not.
    IF NOT FOUND THEN
        RAISE EXCEPTION 'No flight at %.', $1;
    END IF;

    RETURN;
END;
$BODY$
LANGUAGE plpgsql;

-- Returns available flights or raises exception if there are no
-- available flights.
SELECT * FROM get_available_flightid(CURRENT_DATE);
```

Note

The current implementation of RETURN NEXT and RETURN QUERY stores the entire result set before returning from the function, as discussed above. That means that if a PL/pgSQL function produces a very large result set, performance might be poor: data will be written to disk to avoid memory exhaustion, but the function itself will not return until the entire result set has been generated. A future version of PL/pgSQL might allow users to define set-returning functions that do not have this limitation. Currently, the point at which data begins being written to disk is controlled by the [work_mem](#) configuration variable. Administrators who have sufficient memory to store larger result sets in memory should consider increasing this parameter.

44.6.2. Returning from a Procedure

A procedure does not have a return value. A procedure can therefore end without a `RETURN` statement. If you wish to use a `RETURN` statement to exit the code early, write just `RETURN` with no expression.

If the procedure has output parameters, the final values of the output parameter variables will be returned to the caller.

44.6.3. Calling a Procedure

A PL/pgSQL function, procedure, or `DO` block can call a procedure using `CALL`. Output parameters are handled differently from the way that `CALL` works in plain SQL. Each `OUT` or `INOUT` parameter of the procedure must correspond to a variable in the `CALL` statement, and whatever the procedure returns is assigned back to that variable after it returns. For example:

```
CREATE PROCEDURE triple(INOUT x int)
LANGUAGE plpgsql
AS $$
BEGIN
    x := x * 3;
END;
$$;

DO $$
DECLARE myvar int := 5;
BEGIN
    CALL triple(myvar);
    RAISE NOTICE 'myvar = %', myvar; -- prints 15
END;
$$;
```

The variable corresponding to an output parameter can be a simple variable or a field of a composite-type variable. Currently, it cannot be an element of an array.

44.6.4. Conditionals

`IF` and `CASE` statements let you execute alternative commands based on certain conditions. PL/pgSQL has three forms of `IF`:

- `IF ... THEN ... END IF`
- `IF ... THEN ... ELSE ... END IF`
- `IF ... THEN ... ELSIF ... THEN ... ELSE ... END IF`

and two forms of `CASE`:

- `CASE ... WHEN ... THEN ... ELSE ... END CASE`
- `CASE WHEN ... THEN ... ELSE ... END CASE`

44.6.4.1. IF-THEN

```
IF boolean-expression THEN
    statements
END IF;
```

`IF-THEN` statements are the simplest form of `IF`. The statements between `THEN` and `END IF` will be executed if the condition is true. Otherwise, they are skipped.

Example:

```
IF v_user_id <> 0 THEN
```

```
UPDATE users SET email = v_email WHERE user_id = v_user_id;
END IF;
```

44.6.4.2. IF-THEN-ELSE

```
IF boolean-expression THEN
    statements
ELSE
    statements
END IF;
```

IF-THEN-ELSE statements add to IF-THEN by letting you specify an alternative set of statements that should be executed if the condition is not true. (Note this includes the case where the condition evaluates to NULL.)

Examples:

```
IF parentid IS NULL OR parentid = ''
THEN
    RETURN fullname;
ELSE
    RETURN hp_true_filename(parentid) || '/' || fullname;
END IF;

IF v_count > 0 THEN
    INSERT INTO users_count (count) VALUES (v_count);
    RETURN 't';
ELSE
    RETURN 'f';
END IF;
```

44.6.4.3. IF-THEN-ELSIF

```
IF boolean-expression THEN
    statements
[ ELIF boolean-expression THEN
    statements
[ ELIF boolean-expression THEN
    statements
    ...
]
[ ELSE
    statements ]
END IF;
```

Sometimes there are more than just two alternatives. IF-THEN-ELSIF provides a convenient method of checking several alternatives in turn. The IF conditions are tested successively until the first one that is true is found. Then the associated statement(s) are executed, after which control passes to the next statement after END IF. (Any subsequent IF conditions are *not* tested.) If none of the IF conditions is true, then the ELSE block (if any) is executed.

Here is an example:

```
IF number = 0 THEN
    result := 'zero';
ELSIF number > 0 THEN
    result := 'positive';
ELSIF number < 0 THEN
    result := 'negative';
ELSE
```

```
-- hmm, the only other possibility is that number is null
result := 'NULL';
END IF;
```

The key word `ELSIF` can also be spelled `ELSEIF`.

An alternative way of accomplishing the same task is to nest `IF-THEN-ELSE` statements, as in the following example:

```
IF demo_row.sex = 'm' THEN
    pretty_sex := 'man';
ELSE
    IF demo_row.sex = 'f' THEN
        pretty_sex := 'woman';
    END IF;
END IF;
```

However, this method requires writing a matching `END IF` for each `IF`, so it is much more cumbersome than using `ELSIF` when there are many alternatives.

44.6.4.4. Simple CASE

```
CASE search-expression
    WHEN expression [, expression [ ... ]] THEN
        statements
    [ WHEN expression [, expression [ ... ]] THEN
        statements
    ... ]
    [ ELSE
        statements ]
END CASE;
```

The simple form of `CASE` provides conditional execution based on equality of operands. The *search-expression* is evaluated (once) and successively compared to each *expression* in the `WHEN` clauses. If a match is found, then the corresponding *statements* are executed, and then control passes to the next statement after `END CASE`. (Subsequent `WHEN` expressions are not evaluated.) If no match is found, the `ELSE statements` are executed; but if `ELSE` is not present, then a `CASE_NOT_FOUND` exception is raised.

Here is a simple example:

```
CASE x
    WHEN 1, 2 THEN
        msg := 'one or two';
    ELSE
        msg := 'other value than one or two';
END CASE;
```

44.6.4.5. Searched CASE

```
CASE
    WHEN boolean-expression THEN
        statements
    [ WHEN boolean-expression THEN
        statements
    ... ]
    [ ELSE
        statements ]
END CASE;
```

The searched form of `CASE` provides conditional execution based on truth of Boolean expressions. Each `WHEN` clause's *boolean-expression* is evaluated in turn, until one is found that yields `true`. Then the

corresponding *statements* are executed, and then control passes to the next statement after `END CASE`. (Subsequent `WHEN` expressions are not evaluated.) If no true result is found, the `ELSE statements` are executed; but if `ELSE` is not present, then a `CASE_NOT_FOUND` exception is raised.

Here is an example:

```
CASE
  WHEN x BETWEEN 0 AND 10 THEN
    msg := 'value is between zero and ten';
  WHEN x BETWEEN 11 AND 20 THEN
    msg := 'value is between eleven and twenty';
END CASE;
```

This form of `CASE` is entirely equivalent to `IF-THEN-ELSIF`, except for the rule that reaching an omitted `ELSE` clause results in an error rather than doing nothing.

44.6.5. Simple Loops

With the `LOOP`, `EXIT`, `CONTINUE`, `WHILE`, `FOR`, and `FOREACH` statements, you can arrange for your PL/pgSQL function to repeat a series of commands.

44.6.5.1. LOOP

```
[ <<label>> ]
LOOP
  statements
END LOOP [ label ];
```

`LOOP` defines an unconditional loop that is repeated indefinitely until terminated by an `EXIT` or `RETURN` statement. The optional *label* can be used by `EXIT` and `CONTINUE` statements within nested loops to specify which loop those statements refer to.

44.6.5.2. EXIT

```
EXIT [ label ] [ WHEN boolean-expression ];
```

If no *label* is given, the innermost loop is terminated and the statement following `END LOOP` is executed next. If *label* is given, it must be the label of the current or some outer level of nested loop or block. Then the named loop or block is terminated and control continues with the statement after the loop's/block's corresponding `END`.

If `WHEN` is specified, the loop exit occurs only if *boolean-expression* is true. Otherwise, control passes to the statement after `EXIT`.

`EXIT` can be used with all types of loops; it is not limited to use with unconditional loops.

When used with a `BEGIN` block, `EXIT` passes control to the next statement after the end of the block. Note that a label must be used for this purpose; an unlabeled `EXIT` is never considered to match a `BEGIN` block. (This is a change from pre-8.4 releases of PostgreSQL, which would allow an unlabeled `EXIT` to match a `BEGIN` block.)

Examples:

```
LOOP
  -- some computations
  IF count > 0 THEN
    EXIT; -- exit loop
  END IF;
END LOOP;
```

```

LOOP
    -- some computations
    EXIT WHEN count > 0; -- same result as previous example
END LOOP;

<<ablock>>
BEGIN
    -- some computations
    IF stocks > 100000 THEN
        EXIT ablock; -- causes exit from the BEGIN block
    END IF;
    -- computations here will be skipped when stocks > 100000
END;

```

44.6.5.3. CONTINUE

```
CONTINUE [ label ] [ WHEN boolean-expression ];
```

If no *label* is given, the next iteration of the innermost loop is begun. That is, all statements remaining in the loop body are skipped, and control returns to the loop control expression (if any) to determine whether another loop iteration is needed. If *label* is present, it specifies the label of the loop whose execution will be continued.

If **WHEN** is specified, the next iteration of the loop is begun only if *boolean-expression* is true. Otherwise, control passes to the statement after **CONTINUE**.

CONTINUE can be used with all types of loops; it is not limited to use with unconditional loops.

Examples:

```

LOOP
    -- some computations
    EXIT WHEN count > 100;
    CONTINUE WHEN count < 50;
    -- some computations for count IN [50 .. 100]
END LOOP;

```

44.6.5.4. WHILE

```

[ <<label>> ]
WHILE boolean-expression LOOP
    statements
END LOOP [ label ];

```

The **WHILE** statement repeats a sequence of statements so long as the *boolean-expression* evaluates to true. The expression is checked just before each entry to the loop body.

For example:

```

WHILE amount_owed > 0 AND gift_certificate_balance > 0 LOOP
    -- some computations here
END LOOP;

WHILE NOT done LOOP
    -- some computations here
END LOOP;

```

44.6.5.5. FOR (Integer Variant)

```

[ <<label>> ]
FOR name IN [ REVERSE ] expression .. expression [ BY expression ] LOOP

```

```
statements
END LOOP [ label ];
```

This form of `FOR` creates a loop that iterates over a range of integer values. The variable *name* is automatically defined as type `integer` and exists only inside the loop (any existing definition of the variable *name* is ignored within the loop). The two expressions giving the lower and upper bound of the range are evaluated once when entering the loop. If the `BY` clause isn't specified the iteration step is 1, otherwise it's the value specified in the `BY` clause, which again is evaluated once on loop entry. If `REVERSE` is specified then the step value is subtracted, rather than added, after each iteration.

Some examples of integer `FOR` loops:

```
FOR i IN 1..10 LOOP
    -- i will take on the values 1,2,3,4,5,6,7,8,9,10 within the loop
END LOOP;
```

```
FOR i IN REVERSE 10..1 LOOP
    -- i will take on the values 10,9,8,7,6,5,4,3,2,1 within the loop
END LOOP;
```

```
FOR i IN REVERSE 10..1 BY 2 LOOP
    -- i will take on the values 10,8,6,4,2 within the loop
END LOOP;
```

If the lower bound is greater than the upper bound (or less than, in the `REVERSE` case), the loop body is not executed at all. No error is raised.

If a *label* is attached to the `FOR` loop then the integer loop variable can be referenced with a qualified name, using that *label*.

44.6.6. Looping through Query Results

Using a different type of `FOR` loop, you can iterate through the results of a query and manipulate that data accordingly. The syntax is:

```
[ <<label>> ]
FOR target IN query LOOP
    statements
END LOOP [ label ];
```

The *target* is a record variable, row variable, or comma-separated list of scalar variables. The *target* is successively assigned each row resulting from the *query* and the loop body is executed for each row. Here is an example:

```
CREATE FUNCTION refresh_mviews() RETURNS integer AS $$
DECLARE
    mviews RECORD;
BEGIN
    RAISE NOTICE 'Refreshing all materialized views...';

    FOR mviews IN
        SELECT n.nspname AS mv_schema,
               c.relname AS mv_name,
               pg_catalog.pg_get_userbyid(c.relowner) AS owner
        FROM pg_catalog.pg_class c
        LEFT JOIN pg_catalog.pg_namespace n ON (n.oid = c.relnamespace)
        WHERE c.relkind = 'm'
        ORDER BY 1
    LOOP
```

```
-- Now "mviews" has one record with information about the materialized view

RAISE NOTICE 'Refreshing materialized view %.% (owner: %)...',
    quote_ident(mviews.mv_schema),
    quote_ident(mviews.mv_name),
    quote_ident(mviews.owner);
EXECUTE format('REFRESH MATERIALIZED VIEW %I.%I', mviews.mv_schema,
mviews.mv_name);
END LOOP;

RAISE NOTICE 'Done refreshing materialized views.';
RETURN 1;
END;
$$ LANGUAGE plpgsql;
```

If the loop is terminated by an `EXIT` statement, the last assigned row value is still accessible after the loop.

The *query* used in this type of `FOR` statement can be any SQL command that returns rows to the caller: `SELECT` is the most common case, but you can also use `INSERT`, `UPDATE`, `DELETE`, or `MERGE` with a `RETURNING` clause. Some utility commands such as `EXPLAIN` will work too.

PL/pgSQL variables are replaced by query parameters, and the query plan is cached for possible re-use, as discussed in detail in [Section 44.12.1](#) and [Section 44.12.2](#).

The `FOR-IN-EXECUTE` statement is another way to iterate over rows:

```
[ <<label>> ]
FOR target IN EXECUTE text_expression [ USING expression [, ... ] ] LOOP
    statements
END LOOP [ label ];
```

This is like the previous form, except that the source query is specified as a string expression, which is evaluated and replanned on each entry to the `FOR` loop. This allows the programmer to choose the speed of a preplanned query or the flexibility of a dynamic query, just as with a plain `EXECUTE` statement. As with `EXECUTE`, parameter values can be inserted into the dynamic command via `USING`.

Another way to specify the query whose results should be iterated through is to declare it as a cursor. This is described in [Section 44.7.4](#).

44.6.7. Looping through Arrays

The `FOREACH` loop is much like a `FOR` loop, but instead of iterating through the rows returned by an SQL query, it iterates through the elements of an array value. (In general, `FOREACH` is meant for looping through components of a composite-valued expression; variants for looping through composites besides arrays may be added in future.) The `FOREACH` statement to loop over an array is:

```
[ <<label>> ]
FOREACH target [ SLICE number ] IN ARRAY expression LOOP
    statements
END LOOP [ label ];
```

Without `SLICE`, or if `SLICE 0` is specified, the loop iterates through individual elements of the array produced by evaluating the *expression*. The *target* variable is assigned each element value in sequence, and the loop body is executed for each element. Here is an example of looping through the elements of an integer array:

```
CREATE FUNCTION sum(int[]) RETURNS int8 AS $$
DECLARE
    s int8 := 0;
```

```

    x int;
BEGIN
    FOREACH x IN ARRAY $1
    LOOP
        s := s + x;
    END LOOP;
    RETURN s;
END;
$$ LANGUAGE plpgsql;

```

The elements are visited in storage order, regardless of the number of array dimensions. Although the *target* is usually just a single variable, it can be a list of variables when looping through an array of composite values (records). In that case, for each array element, the variables are assigned from successive columns of the composite value.

With a positive *SLICE* value, *FOREACH* iterates through slices of the array rather than single elements. The *SLICE* value must be an integer constant not larger than the number of dimensions of the array. The *target* variable must be an array, and it receives successive slices of the array value, where each slice is of the number of dimensions specified by *SLICE*. Here is an example of iterating through one-dimensional slices:

```

CREATE FUNCTION scan_rows(int[]) RETURNS void AS $$
DECLARE
    x int[];
BEGIN
    FOREACH x SLICE 1 IN ARRAY $1
    LOOP
        RAISE NOTICE 'row = %', x;
    END LOOP;
END;
$$ LANGUAGE plpgsql;

SELECT scan_rows(ARRAY[[1,2,3],[4,5,6],[7,8,9],[10,11,12]]);

NOTICE:  row = {1,2,3}
NOTICE:  row = {4,5,6}
NOTICE:  row = {7,8,9}
NOTICE:  row = {10,11,12}

```

44.6.8. Trapping Errors

By default, any error occurring in a PL/pgSQL function aborts execution of the function and the surrounding transaction. You can trap errors and recover from them by using a *BEGIN* block with an *EXCEPTION* clause. The syntax is an extension of the normal syntax for a *BEGIN* block:

```

[ <<label>> ]
[ DECLARE
    declarations ]
BEGIN
    statements
EXCEPTION
    WHEN condition [ OR condition ... ] THEN
        handler_statements
    [ WHEN condition [ OR condition ... ] THEN
        handler_statements
    ... ]
END;

```

If no error occurs, this form of block simply executes all the *statements*, and then control passes to the next statement after *END*. But if an error occurs within the *statements*, further processing of the

statements is abandoned, and control passes to the `EXCEPTION` list. The list is searched for the first *condition* matching the error that occurred. If a match is found, the corresponding *handler_statements* are executed, and then control passes to the next statement after `END`. If no match is found, the error propagates out as though the `EXCEPTION` clause were not there at all: the error can be caught by an enclosing block with `EXCEPTION`, or if there is none it aborts processing of the function.

The *condition* names can be any of those shown in [Appendix A](#). A category name matches any error within its category. The special condition name `OTHERS` matches every error type except `QUERY_CANCELED` and `ASSERT_FAILURE`. (It is possible, but often unwise, to trap those two error types by name.) Condition names are not case-sensitive. Also, an error condition can be specified by `SQLSTATE` code; for example these are equivalent:

```
WHEN division_by_zero THEN ...
WHEN SQLSTATE '22012' THEN ...
```

If a new error occurs within the selected *handler_statements*, it cannot be caught by this `EXCEPTION` clause, but is propagated out. A surrounding `EXCEPTION` clause could catch it.

When an error is caught by an `EXCEPTION` clause, the local variables of the PL/pgSQL function remain as they were when the error occurred, but all changes to persistent database state within the block are rolled back. As an example, consider this fragment:

```
INSERT INTO mytab(firstname, lastname) VALUES('Tom', 'Jones');
BEGIN
    UPDATE mytab SET firstname = 'Joe' WHERE lastname = 'Jones';
    x := x + 1;
    y := x / 0;
EXCEPTION
    WHEN division_by_zero THEN
        RAISE NOTICE 'caught division_by_zero';
        RETURN x;
END;
```

When control reaches the assignment to `y`, it will fail with a `division_by_zero` error. This will be caught by the `EXCEPTION` clause. The value returned in the `RETURN` statement will be the incremented value of `x`, but the effects of the `UPDATE` command will have been rolled back. The `INSERT` command preceding the block is not rolled back, however, so the end result is that the database contains Tom Jones not Joe Jones.

Tip

A block containing an `EXCEPTION` clause is significantly more expensive to enter and exit than a block without one. Therefore, don't use `EXCEPTION` without need.

Example 44.2. Exceptions with `UPDATE/INSERT`

This example uses exception handling to perform either `UPDATE` or `INSERT`, as appropriate. It is recommended that applications use `INSERT` with `ON CONFLICT DO UPDATE` rather than actually using this pattern. This example serves primarily to illustrate use of PL/pgSQL control flow structures:

```
CREATE TABLE db (a INT PRIMARY KEY, b TEXT);

CREATE FUNCTION merge_db(key INT, data TEXT) RETURNS VOID AS
$$
BEGIN
    LOOP
        -- first try to update the key
```

```

UPDATE db SET b = data WHERE a = key;
IF found THEN
    RETURN;
END IF;
-- not there, so try to insert the key
-- if someone else inserts the same key concurrently,
-- we could get a unique-key failure
BEGIN
    INSERT INTO db(a,b) VALUES (key, data);
    RETURN;
EXCEPTION WHEN unique_violation THEN
    -- Do nothing, and loop to try the UPDATE again.
END;
END LOOP;
END;
$$
LANGUAGE plpgsql;

SELECT merge_db(1, 'david');
SELECT merge_db(1, 'dennis');
```

This coding assumes the `unique_violation` error is caused by the `INSERT`, and not by, say, an `INSERT` in a trigger function on the table. It might also misbehave if there is more than one unique index on the table, since it will retry the operation regardless of which index caused the error. More safety could be had by using the features discussed next to check that the trapped error was the one expected.

44.6.8.1. Obtaining Information about an Error

Exception handlers frequently need to identify the specific error that occurred. There are two ways to get information about the current exception in PL/pgSQL: special variables and the `GET STACKED DIAGNOSTICS` command.

Within an exception handler, the special variable `SQLSTATE` contains the error code that corresponds to the exception that was raised (refer to [Table A.1](#) for a list of possible error codes). The special variable `SQLERRM` contains the error message associated with the exception. These variables are undefined outside exception handlers.

Within an exception handler, one may also retrieve information about the current exception by using the `GET STACKED DIAGNOSTICS` command, which has the form:

```
GET STACKED DIAGNOSTICS variable { = | := } item [ , ... ];
```

Each *item* is a key word identifying a status value to be assigned to the specified *variable* (which should be of the right data type to receive it). The currently available status items are shown in [Table 44.2](#).

Table 44.2. Error Diagnostics Items

Name	Type	Description
RETURNED_SQLSTATE	text	the SQLSTATE error code of the exception
COLUMN_NAME	text	the name of the column related to exception
CONSTRAINT_NAME	text	the name of the constraint related to exception
PG_DATATYPE_NAME	text	the name of the data type related to exception
MESSAGE_TEXT	text	the text of the exception's primary message

Name	Type	Description
TABLE_NAME	text	the name of the table related to exception
SCHEMA_NAME	text	the name of the schema related to exception
PG_EXCEPTION_DETAIL	text	the text of the exception's detail message, if any
PG_EXCEPTION_HINT	text	the text of the exception's hint message, if any
PG_EXCEPTION_CONTEXT	text	line(s) of text describing the call stack at the time of the exception (see Section 44.6.9)

If the exception did not set a value for an item, an empty string will be returned.

Here is an example:

```
DECLARE
    text_var1 text;
    text_var2 text;
    text_var3 text;
BEGIN
    -- some processing which might cause an exception
    ...
EXCEPTION WHEN OTHERS THEN
    GET STACKED DIAGNOSTICS text_var1 = MESSAGE_TEXT,
                           text_var2 = PG_EXCEPTION_DETAIL,
                           text_var3 = PG_EXCEPTION_HINT;
END;
```

44.6.9. Obtaining Execution Location Information

The `GET DIAGNOSTICS` command, previously described in [Section 44.5.5](#), retrieves information about current execution state (whereas the `GET STACKED DIAGNOSTICS` command discussed above reports information about the execution state as of a previous error). Its `PG_CONTEXT` status item is useful for identifying the current execution location. `PG_CONTEXT` returns a text string with line(s) of text describing the call stack. The first line refers to the current function and currently executing `GET DIAGNOSTICS` command. The second and any subsequent lines refer to calling functions further up the call stack. For example:

```
CREATE OR REPLACE FUNCTION outer_func() RETURNS integer AS $$
BEGIN
    RETURN inner_func();
END;
$$ LANGUAGE plpgsql;

CREATE OR REPLACE FUNCTION inner_func() RETURNS integer AS $$
DECLARE
    stack text;
BEGIN
    GET DIAGNOSTICS stack = PG_CONTEXT;
    RAISE NOTICE E'--- Call Stack ---\n%', stack;
    RETURN 1;
END;
$$ LANGUAGE plpgsql;
```

```
SELECT outer_func();

NOTICE: --- Call Stack ---
PL/pgSQL function inner_func() line 5 at GET DIAGNOSTICS
PL/pgSQL function outer_func() line 3 at RETURN
CONTEXT: PL/pgSQL function outer_func() line 3 at RETURN
outer_func
-----
          1
(1 row)
```

GET STACKED DIAGNOSTICS ... PG_EXCEPTION_CONTEXT returns the same sort of stack trace, but describing the location at which an error was detected, rather than the current location.

44.7. Cursors

Rather than executing a whole query at once, it is possible to set up a *cursor* that encapsulates the query, and then read the query result a few rows at a time. One reason for doing this is to avoid memory overrun when the result contains a large number of rows. (However, PL/pgSQL users do not normally need to worry about that, since FOR loops automatically use a cursor internally to avoid memory problems.) A more interesting usage is to return a reference to a cursor that a function has created, allowing the caller to read the rows. This provides an efficient way to return large row sets from functions.

44.7.1. Declaring Cursor Variables

All access to cursors in PL/pgSQL goes through cursor variables, which are always of the special data type `refcursor`. One way to create a cursor variable is just to declare it as a variable of type `refcursor`. Another way is to use the cursor declaration syntax, which in general is:

```
name [ [ NO ] SCROLL ] CURSOR [ ( arguments ) ] FOR query;
```

(FOR can be replaced by IS for Oracle compatibility.) If SCROLL is specified, the cursor will be capable of scrolling backward; if NO SCROLL is specified, backward fetches will be rejected; if neither specification appears, it is query-dependent whether backward fetches will be allowed. *arguments*, if specified, is a comma-separated list of pairs *name datatype* that define names to be replaced by parameter values in the given query. The actual values to substitute for these names will be specified later, when the cursor is opened.

Some examples:

```
DECLARE
    curs1 refcursor;
    curs2 CURSOR FOR SELECT * FROM tenk1;
    curs3 CURSOR (key integer) FOR SELECT * FROM tenk1 WHERE unique1 = key;
```

All three of these variables have the data type `refcursor`, but the first can be used with any query, while the second has a fully specified query already *bound* to it, and the last has a parameterized query bound to it. (*key* will be replaced by an integer parameter value when the cursor is opened.) The variable `curs1` is said to be *unbound* since it is not bound to any particular query.

The SCROLL option cannot be used when the cursor's query uses FOR UPDATE/SHARE. Also, it is best to use NO SCROLL with a query that involves volatile functions. The implementation of SCROLL assumes that re-reading the query's output will give consistent results, which a volatile function might not do.

44.7.2. Opening Cursors

Before a cursor can be used to retrieve rows, it must be *opened*. (This is the equivalent action to the SQL command `DECLARE CURSOR`.) PL/pgSQL has three forms of the OPEN statement, two of which use unbound cursor variables while the third uses a bound cursor variable.

Note

Bound cursor variables can also be used without explicitly opening the cursor, via the `FOR` statement described in [Section 44.7.4](#). A `FOR` loop will open the cursor and then close it again when the loop completes.

Opening a cursor involves creating a server-internal data structure called a *portal*, which holds the execution state for the cursor's query. A portal has a name, which must be unique within the session for the duration of the portal's existence. By default, PL/pgSQL will assign a unique name to each portal it creates. However, if you assign a non-null string value to a cursor variable, that string will be used as its portal name. This feature can be used as described in [Section 44.7.3.5](#).

44.7.2.1. OPEN FOR query

```
OPEN unbound_cursorvar [ [ NO ] SCROLL ] FOR query;
```

The cursor variable is opened and given the specified query to execute. The cursor cannot be open already, and it must have been declared as an unbound cursor variable (that is, as a simple `refcursor` variable). The query must be a `SELECT`, or something else that returns rows (such as `EXPLAIN`). The query is treated in the same way as other SQL commands in PL/pgSQL: PL/pgSQL variable names are substituted, and the query plan is cached for possible reuse. When a PL/pgSQL variable is substituted into the cursor query, the value that is substituted is the one it has at the time of the `OPEN`; subsequent changes to the variable will not affect the cursor's behavior. The `SCROLL` and `NO SCROLL` options have the same meanings as for a bound cursor.

An example:

```
OPEN curs1 FOR SELECT * FROM foo WHERE key = mykey;
```

44.7.2.2. OPEN FOR EXECUTE

```
OPEN unbound_cursorvar [ [ NO ] SCROLL ] FOR EXECUTE query_string
                        [ USING expression [, ... ] ];
```

The cursor variable is opened and given the specified query to execute. The cursor cannot be open already, and it must have been declared as an unbound cursor variable (that is, as a simple `refcursor` variable). The query is specified as a string expression, in the same way as in the `EXECUTE` command. As usual, this gives flexibility so the query plan can vary from one run to the next (see [Section 44.12.2](#)), and it also means that variable substitution is not done on the command string. As with `EXECUTE`, parameter values can be inserted into the dynamic command via `format()` and `USING`. The `SCROLL` and `NO SCROLL` options have the same meanings as for a bound cursor.

An example:

```
OPEN curs1 FOR EXECUTE format('SELECT * FROM %I WHERE col1 = $1',tablename) USING
keyvalue;
```

In this example, the table name is inserted into the query via `format()`. The comparison value for `col1` is inserted via a `USING` parameter, so it needs no quoting.

44.7.2.3. Opening a Bound Cursor

```
OPEN bound_cursorvar [ ( [ argument_name := ] argument_value [, ... ] ) ];
```

This form of `OPEN` is used to open a cursor variable whose query was bound to it when it was declared. The cursor cannot be open already. A list of actual argument value expressions must appear if and only if the cursor was declared to take arguments. These values will be substituted in the query.

The query plan for a bound cursor is always considered cacheable; there is no equivalent of `EXECUTE` in this case. Notice that `SCROLL` and `NO SCROLL` cannot be specified in `OPEN`, as the cursor's scrolling behavior was already determined.

Argument values can be passed using either *positional* or *named* notation. In positional notation, all arguments are specified in order. In named notation, each argument's name is specified using `:=` to separate it from the argument expression. Similar to calling functions, described in [Section 4.3](#), it is also allowed to mix positional and named notation.

Examples (these use the cursor declaration examples above):

```
OPEN curs2;
OPEN curs3(42);
OPEN curs3(key := 42);
```

Because variable substitution is done on a bound cursor's query, there are really two ways to pass values into the cursor: either with an explicit argument to `OPEN`, or implicitly by referencing a PL/pgSQL variable in the query. However, only variables declared before the bound cursor was declared will be substituted into it. In either case the value to be passed is determined at the time of the `OPEN`. For example, another way to get the same effect as the `curs3` example above is

```
DECLARE
    key integer;
    curs4 CURSOR FOR SELECT * FROM tenk1 WHERE unique1 = key;
BEGIN
    key := 42;
    OPEN curs4;
```

44.7.3. Using Cursors

Once a cursor has been opened, it can be manipulated with the statements described here.

These manipulations need not occur in the same function that opened the cursor to begin with. You can return a `refcursor` value out of a function and let the caller operate on the cursor. (Internally, a `refcursor` value is simply the string name of the portal containing the active query for the cursor. This name can be passed around, assigned to other `refcursor` variables, and so on, without disturbing the portal.)

All portals are implicitly closed at transaction end. Therefore a `refcursor` value is usable to reference an open cursor only until the end of the transaction.

44.7.3.1. FETCH

```
FETCH [ direction { FROM | IN } ] cursor INTO target;
```

`FETCH` retrieves the next row (in the indicated direction) from the cursor into a target, which might be a row variable, a record variable, or a comma-separated list of simple variables, just like `SELECT INTO`. If there is no suitable row, the target is set to `NULL(s)`. As with `SELECT INTO`, the special variable `FOUND` can be checked to see whether a row was obtained or not. If no row is obtained, the cursor is positioned after the last row or before the first row, depending on the movement direction.

The *direction* clause can be any of the variants allowed in the SQL [FETCH](#) command except the ones that can fetch more than one row; namely, it can be `NEXT`, `PRIOR`, `FIRST`, `LAST`, `ABSOLUTE count`, `RELATIVE count`, `FORWARD`, or `BACKWARD`. Omitting *direction* is the same as specifying `NEXT`. In the forms using a *count*, the *count* can be any integer-valued expression (unlike the SQL `FETCH` command, which only allows an integer constant). *direction* values that require moving backward are likely to fail unless the cursor was declared or opened with the `SCROLL` option.

cursor must be the name of a `refcursor` variable that references an open cursor portal.

Examples:

```
FETCH curs1 INTO rowvar;
FETCH curs2 INTO foo, bar, baz;
```

```
FETCH LAST FROM curs3 INTO x, y;
FETCH RELATIVE -2 FROM curs4 INTO x;
```

44.7.3.2. MOVE

```
MOVE [ direction { FROM | IN } ] cursor;
```

`MOVE` repositions a cursor without retrieving any data. `MOVE` works like the `FETCH` command, except it only repositions the cursor and does not return the row moved to. The *direction* clause can be any of the variants allowed in the SQL `FETCH` command, including those that can fetch more than one row; the cursor is positioned to the last such row. (However, the case in which the *direction* clause is simply a *count* expression with no key word is deprecated in PL/pgSQL. That syntax is ambiguous with the case where the *direction* clause is omitted altogether, and hence it may fail if the *count* is not a constant.) As with `SELECT INTO`, the special variable `FOUND` can be checked to see whether there was a row to move to. If there is no such row, the cursor is positioned after the last row or before the first row, depending on the movement direction.

Examples:

```
MOVE curs1;
MOVE LAST FROM curs3;
MOVE RELATIVE -2 FROM curs4;
MOVE FORWARD 2 FROM curs4;
```

44.7.3.3. UPDATE/DELETE WHERE CURRENT OF

```
UPDATE table SET ... WHERE CURRENT OF cursor;
DELETE FROM table WHERE CURRENT OF cursor;
```

When a cursor is positioned on a table row, that row can be updated or deleted using the cursor to identify the row. There are restrictions on what the cursor's query can be (in particular, no grouping) and it's best to use `FOR UPDATE` in the cursor. For more information see the [DECLARE](#) reference page.

An example:

```
UPDATE foo SET dataval = myval WHERE CURRENT OF curs1;
```

44.7.3.4. CLOSE

```
CLOSE cursor;
```

`CLOSE` closes the portal underlying an open cursor. This can be used to release resources earlier than end of transaction, or to free up the cursor variable to be opened again.

An example:

```
CLOSE curs1;
```

44.7.3.5. Returning Cursors

PL/pgSQL functions can return cursors to the caller. This is useful to return multiple rows or columns, especially with very large result sets. To do this, the function opens the cursor and returns the cursor name to the caller (or simply opens the cursor using a portal name specified by or otherwise known to the caller). The caller can then fetch rows from the cursor. The cursor can be closed by the caller, or it will be closed automatically when the transaction closes.

The portal name used for a cursor can be specified by the programmer or automatically generated. To specify a portal name, simply assign a string to the `refcursor` variable before opening it. The string value of the `refcursor` variable will be used by `OPEN` as the name of the underlying portal. However, if the `refcursor` variable's value is null (as it will be by default), then `OPEN` automatically generates a name that does not conflict with any existing portal, and assigns it to the `refcursor` variable.

Note

Prior to Postgres Pro 16, bound cursor variables were initialized to contain their own names, rather than being left as null, so that the underlying portal name would be the same as the cursor variable's name by default. This was changed because it created too much risk of conflicts between similarly-named cursors in different functions.

The following example shows one way a cursor name can be supplied by the caller:

```
CREATE TABLE test (col text);
INSERT INTO test VALUES ('123');

CREATE FUNCTION reffunc(refcursor) RETURNS refcursor AS '
BEGIN
    OPEN $1 FOR SELECT col FROM test;
    RETURN $1;
END;
' LANGUAGE plpgsql;

BEGIN;
SELECT reffunc('funccursor');
FETCH ALL IN funccursor;
COMMIT;
```

The following example uses automatic cursor name generation:

```
CREATE FUNCTION reffunc2() RETURNS refcursor AS '
DECLARE
    ref refcursor;
BEGIN
    OPEN ref FOR SELECT col FROM test;
    RETURN ref;
END;
' LANGUAGE plpgsql;

-- need to be in a transaction to use cursors.
BEGIN;
SELECT reffunc2();

      reffunc2
-----
<unnamed cursor 1>
(1 row)

FETCH ALL IN "<unnamed cursor 1>";
COMMIT;
```

The following example shows one way to return multiple cursors from a single function:

```
CREATE FUNCTION myfunc(refcursor, refcursor) RETURNS SETOF refcursor AS $$
BEGIN
    OPEN $1 FOR SELECT * FROM table_1;
    RETURN NEXT $1;
    OPEN $2 FOR SELECT * FROM table_2;
    RETURN NEXT $2;
END;
$$ LANGUAGE plpgsql;
```



```
-- need to be in a transaction to use cursors.
BEGIN;

SELECT * FROM myfunc('a', 'b');

FETCH ALL FROM a;
FETCH ALL FROM b;
COMMIT;
```

44.7.4. Looping through a Cursor's Result

There is a variant of the `FOR` statement that allows iterating through the rows returned by a cursor. The syntax is:

```
[ <<label>> ]
FOR recordvar IN bound_cursorvar [ ( [ argument_name := ] argument_value [, ...] ) ]
    LOOP
        statements
    END LOOP [ label ];
```

The cursor variable must have been bound to some query when it was declared, and it *cannot* be open already. The `FOR` statement automatically opens the cursor, and it closes the cursor again when the loop exits. A list of actual argument value expressions must appear if and only if the cursor was declared to take arguments. These values will be substituted in the query, in just the same way as during an `OPEN` (see [Section 44.7.2.3](#)).

The variable `recordvar` is automatically defined as type `record` and exists only inside the loop (any existing definition of the variable name is ignored within the loop). Each row returned by the cursor is successively assigned to this record variable and the loop body is executed.

44.8. Transaction Management

In procedures invoked by the `CALL` command as well as in anonymous code blocks (`DO` command), it is possible to end transactions using the commands `COMMIT` and `ROLLBACK`. A new transaction is started automatically after a transaction is ended using these commands, so there is no separate `START TRANSACTION` command. (Note that `BEGIN` and `END` have different meanings in PL/pgSQL.)

Here is a simple example:

```
CREATE PROCEDURE transaction_test1()
LANGUAGE plpgsql
AS $$
BEGIN
    FOR i IN 0..9 LOOP
        INSERT INTO test1 (a) VALUES (i);
        IF i % 2 = 0 THEN
            COMMIT;
        ELSE
            ROLLBACK;
        END IF;
    END LOOP;
END;
$$;

CALL transaction_test1();
```

A new transaction starts out with default transaction characteristics such as transaction isolation level. In cases where transactions are committed in a loop, it might be desirable to start new transactions

automatically with the same characteristics as the previous one. The commands `COMMIT AND CHAIN` and `ROLLBACK AND CHAIN` accomplish this.

Transaction control is only possible in `CALL` or `DO` invocations from the top level or nested `CALL` or `DO` invocations without any other intervening command. For example, if the call stack is `CALL proc1() → CALL proc2() → CALL proc3()`, then the second and third procedures can perform transaction control actions. But if the call stack is `CALL proc1() → SELECT func2() → CALL proc3()`, then the last procedure cannot do transaction control, because of the `SELECT` in between.

PL/pgSQL does not support savepoints (`SAVEPOINT/ROLLBACK TO SAVEPOINT/RELEASE SAVEPOINT` commands). Typical usage patterns for savepoints can be replaced by blocks with exception handlers (see [Section 44.6.8](#)). Under the hood, a block with exception handlers forms a subtransaction, which means that transactions cannot be ended inside such a block.

Special considerations apply to cursor loops. Consider this example:

```
CREATE PROCEDURE transaction_test2()
LANGUAGE plpgsql
AS $$
DECLARE
    r RECORD;
BEGIN
    FOR r IN SELECT * FROM test2 ORDER BY x LOOP
        INSERT INTO test1 (a) VALUES (r.x);
        COMMIT;
    END LOOP;
END;
$$;

CALL transaction_test2();
```

Normally, cursors are automatically closed at transaction commit. However, a cursor created as part of a loop like this is automatically converted to a holdable cursor by the first `COMMIT` or `ROLLBACK`. That means that the cursor is fully evaluated at the first `COMMIT` or `ROLLBACK` rather than row by row. The cursor is still removed automatically after the loop, so this is mostly invisible to the user. But one must keep in mind that any table or row locks taken by the cursor's query will no longer be held after the first `COMMIT` or `ROLLBACK`.

Transaction commands are not allowed in cursor loops driven by commands that are not read-only (for example `UPDATE ... RETURNING`).

44.9. Errors and Messages

44.9.1. Reporting Errors and Messages

Use the `RAISE` statement to report messages and raise errors.

```
RAISE [ level ] 'format' [, expression [, ... ] ] [ USING option = expression
[, ... ] ];
RAISE [ level ] condition_name [ USING option = expression [, ... ] ];
RAISE [ level ] SQLSTATE 'sqlstate' [ USING option = expression [, ... ] ];
RAISE [ level ] USING option = expression [, ... ];
RAISE ;
```

The *level* option specifies the error severity. Allowed levels are `DEBUG`, `LOG`, `INFO`, `NOTICE`, `WARNING`, and `EXCEPTION`, with `EXCEPTION` being the default. `EXCEPTION` raises an error (which normally aborts the current transaction); the other levels only generate messages of different priority levels. Whether messages of a particular priority are reported to the client, written to the server log, or both is controlled by the [log_min_messages](#) and [client_min_messages](#) configuration variables. See [Chapter 19](#) for more information.

After *level* if any, you can specify a *format* string (which must be a simple string literal, not an expression). The format string specifies the error message text to be reported. The format string can be followed by optional argument expressions to be inserted into the message. Inside the format string, % is replaced by the string representation of the next optional argument's value. Write %% to emit a literal %. The number of arguments must match the number of % placeholders in the format string, or an error is raised during the compilation of the function.

In this example, the value of `v_job_id` will replace the % in the string:

```
RAISE NOTICE 'Calling cs_create_job(%)', v_job_id;
```

You can attach additional information to the error report by writing `USING` followed by *option = expression* items. Each *expression* can be any string-valued expression. The allowed *option* key words are:

`MESSAGE`

Sets the error message text. This option can't be used in the form of `RAISE` that includes a format string before `USING`.

`DETAIL`

Supplies an error detail message.

`HINT`

Supplies a hint message.

`ERRCODE`

Specifies the error code (SQLSTATE) to report, either by condition name, as shown in [Appendix A](#), or directly as a five-character SQLSTATE code.

`COLUMN`

`CONSTRAINT`

`DATATYPE`

`TABLE`

`SCHEMA`

Supplies the name of a related object.

This example will abort the transaction with the given error message and hint:

```
RAISE EXCEPTION 'Nonexistent ID --> %', user_id
    USING HINT = 'Please check your user ID';
```

These two examples show equivalent ways of setting the SQLSTATE:

```
RAISE 'Duplicate user ID: %', user_id USING ERRCODE = 'unique_violation';
RAISE 'Duplicate user ID: %', user_id USING ERRCODE = '23505';
```

There is a second `RAISE` syntax in which the main argument is the condition name or SQLSTATE to be reported, for example:

```
RAISE division_by_zero;
RAISE SQLSTATE '22012';
```

In this syntax, `USING` can be used to supply a custom error message, detail, or hint. Another way to do the earlier example is

```
RAISE unique_violation USING MESSAGE = 'Duplicate user ID: ' || user_id;
```

Still another variant is to write `RAISE USING` or `RAISE level USING` and put everything else into the `USING` list.

The last variant of `RAISE` has no parameters at all. This form can only be used inside a `BEGIN` block's `EXCEPTION` clause; it causes the error currently being handled to be re-thrown.

Note

Before PostgreSQL 9.1, `RAISE` without parameters was interpreted as re-throwing the error from the block containing the active exception handler. Thus an `EXCEPTION` clause nested within that handler could not catch it, even if the `RAISE` was within the nested `EXCEPTION` clause's block. This was deemed surprising as well as being incompatible with Oracle's PL/SQL.

If no condition name nor `SQLSTATE` is specified in a `RAISE EXCEPTION` command, the default is to use `raise_exception` (P0001). If no message text is specified, the default is to use the condition name or `SQLSTATE` as message text.

Note

When specifying an error code by `SQLSTATE` code, you are not limited to the predefined error codes, but can select any error code consisting of five digits and/or upper-case ASCII letters, other than `00000`. It is recommended that you avoid throwing error codes that end in three zeroes, because these are category codes and can only be trapped by trapping the whole category.

44.9.2. Checking Assertions

The `ASSERT` statement is a convenient shorthand for inserting debugging checks into PL/pgSQL functions.

```
ASSERT condition [ , message ];
```

The *condition* is a Boolean expression that is expected to always evaluate to true; if it does, the `ASSERT` statement does nothing further. If the result is false or null, then an `ASSERT_FAILURE` exception is raised. (If an error occurs while evaluating the *condition*, it is reported as a normal error.)

If the optional *message* is provided, it is an expression whose result (if not null) replaces the default error message text “assertion failed”, should the *condition* fail. The *message* expression is not evaluated in the normal case where the assertion succeeds.

Testing of assertions can be enabled or disabled via the configuration parameter `plpgsql.check_asserts`, which takes a Boolean value; the default is `on`. If this parameter is `off` then `ASSERT` statements do nothing.

Note that `ASSERT` is meant for detecting program bugs, not for reporting ordinary error conditions. Use the `RAISE` statement, described above, for that.

44.10. Trigger Functions

PL/pgSQL can be used to define trigger functions on data changes or database events. A trigger function is created with the `CREATE FUNCTION` command, declaring it as a function with no arguments and a return type of `trigger` (for data change triggers) or `event_trigger` (for database event triggers). Special local variables named `TG_something` are automatically defined to describe the condition that triggered the call.

44.10.1. Triggers on Data Changes

A **data change trigger** is declared as a function with no arguments and a return type of `trigger`. Note that the function must be declared with no arguments even if it expects to receive some arguments specified in `CREATE TRIGGER` — such arguments are passed via `TG_ARGV`, as described below.

When a PL/pgSQL function is called as a trigger, several special variables are created automatically in the top-level block. They are:

`NEW` record

new database row for `INSERT/UPDATE` operations in row-level triggers. This variable is null in statement-level triggers and for `DELETE` operations.

`OLD` record

old database row for `UPDATE/DELETE` operations in row-level triggers. This variable is null in statement-level triggers and for `INSERT` operations.

`TG_NAME` name

name of the trigger which fired.

`TG_WHEN` text

`BEFORE`, `AFTER`, or `INSTEAD OF`, depending on the trigger's definition.

`TG_LEVEL` text

`ROW` or `STATEMENT`, depending on the trigger's definition.

`TG_OP` text

operation for which the trigger was fired: `INSERT`, `UPDATE`, `DELETE`, or `TRUNCATE`.

`TG_RELID` oid (references `pg_class.oid`)

object ID of the table that caused the trigger invocation.

`TG_RELNAME` name

table that caused the trigger invocation. This is now deprecated, and could disappear in a future release. Use `TG_TABLE_NAME` instead.

`TG_TABLE_NAME` name

table that caused the trigger invocation.

`TG_TABLE_SCHEMA` name

schema of the table that caused the trigger invocation.

`TG_NARGS` integer

number of arguments given to the trigger function in the `CREATE TRIGGER` statement.

`TG_ARGV` text []

arguments from the `CREATE TRIGGER` statement. The index counts from 0. Invalid indexes (less than 0 or greater than or equal to `tg_nargs`) result in a null value.

A trigger function must return either `NULL` or a record/row value having exactly the structure of the table the trigger was fired for.

Row-level triggers fired `BEFORE` can return null to signal the trigger manager to skip the rest of the operation for this row (i.e., subsequent triggers are not fired, and the `INSERT/UPDATE/DELETE` does not occur for this row). If a nonnull value is returned then the operation proceeds with that row value. Returning a row value different from the original value of `NEW` alters the row that will be inserted or

updated. Thus, if the trigger function wants the triggering action to succeed normally without altering the row value, `NEW` (or a value equal thereto) has to be returned. To alter the row to be stored, it is possible to replace single values directly in `NEW` and return the modified `NEW`, or to build a complete new record/row to return. In the case of a before-trigger on `DELETE`, the returned value has no direct effect, but it has to be nonnull to allow the trigger action to proceed. Note that `NEW` is null in `DELETE` triggers, so returning that is usually not sensible. The usual idiom in `DELETE` triggers is to return `OLD`.

`INSTEAD OF` triggers (which are always row-level triggers, and may only be used on views) can return null to signal that they did not perform any updates, and that the rest of the operation for this row should be skipped (i.e., subsequent triggers are not fired, and the row is not counted in the rows-affected status for the surrounding `INSERT/UPDATE/DELETE`). Otherwise a nonnull value should be returned, to signal that the trigger performed the requested operation. For `INSERT` and `UPDATE` operations, the return value should be `NEW`, which the trigger function may modify to support `INSERT RETURNING` and `UPDATE RETURNING` (this will also affect the row value passed to any subsequent triggers, or passed to a special `EXCLUDED` alias reference within an `INSERT` statement with an `ON CONFLICT DO UPDATE` clause). For `DELETE` operations, the return value should be `OLD`.

The return value of a row-level trigger fired `AFTER` or a statement-level trigger fired `BEFORE` or `AFTER` is always ignored; it might as well be null. However, any of these types of triggers might still abort the entire operation by raising an error.

[Example 44.3](#) shows an example of a trigger function in PL/pgSQL.

Example 44.3. A PL/pgSQL Trigger Function

This example trigger ensures that any time a row is inserted or updated in the table, the current user name and time are stamped into the row. And it checks that an employee's name is given and that the salary is a positive value.

```
CREATE TABLE emp (
    empname      text,
    salary       integer,
    last_date    timestamp,
    last_user    text
);

CREATE FUNCTION emp_stamp() RETURNS trigger AS $emp_stamp$
BEGIN
    -- Check that empname and salary are given
    IF NEW.empname IS NULL THEN
        RAISE EXCEPTION 'empname cannot be null';
    END IF;
    IF NEW.salary IS NULL THEN
        RAISE EXCEPTION '% cannot have null salary', NEW.empname;
    END IF;

    -- Who works for us when they must pay for it?
    IF NEW.salary < 0 THEN
        RAISE EXCEPTION '% cannot have a negative salary', NEW.empname;
    END IF;

    -- Remember who changed the payroll when
    NEW.last_date := current_timestamp;
    NEW.last_user := current_user;
    RETURN NEW;
END;
$emp_stamp$ LANGUAGE plpgsql;

CREATE TRIGGER emp_stamp BEFORE INSERT OR UPDATE ON emp
```

```
FOR EACH ROW EXECUTE FUNCTION emp_stamp();
```

Another way to log changes to a table involves creating a new table that holds a row for each insert, update, or delete that occurs. This approach can be thought of as auditing changes to a table. [Example 44.4](#) shows an example of an audit trigger function in PL/pgSQL.

Example 44.4. A PL/pgSQL Trigger Function for Auditing

This example trigger ensures that any insert, update or delete of a row in the `emp` table is recorded (i.e., audited) in the `emp_audit` table. The current time and user name are stamped into the row, together with the type of operation performed on it.

```
CREATE TABLE emp (
    empname      text NOT NULL,
    salary       integer
);

CREATE TABLE emp_audit(
    operation     char(1) NOT NULL,
    stamp         timestamp NOT NULL,
    userid        text NOT NULL,
    empname       text NOT NULL,
    salary        integer
);

CREATE OR REPLACE FUNCTION process_emp_audit() RETURNS TRIGGER AS $emp_audit$
BEGIN
    --
    -- Create a row in emp_audit to reflect the operation performed on emp,
    -- making use of the special variable TG_OP to work out the operation.
    --
    IF (TG_OP = 'DELETE') THEN
        INSERT INTO emp_audit SELECT 'D', now(), current_user, OLD.*;
    ELSIF (TG_OP = 'UPDATE') THEN
        INSERT INTO emp_audit SELECT 'U', now(), current_user, NEW.*;
    ELSIF (TG_OP = 'INSERT') THEN
        INSERT INTO emp_audit SELECT 'I', now(), current_user, NEW.*;
    END IF;
    RETURN NULL; -- result is ignored since this is an AFTER trigger
END;
$emp_audit$ LANGUAGE plpgsql;

CREATE TRIGGER emp_audit
AFTER INSERT OR UPDATE OR DELETE ON emp
FOR EACH ROW EXECUTE FUNCTION process_emp_audit();
```

A variation of the previous example uses a view joining the main table to the audit table, to show when each entry was last modified. This approach still records the full audit trail of changes to the table, but also presents a simplified view of the audit trail, showing just the last modified timestamp derived from the audit trail for each entry. [Example 44.5](#) shows an example of an audit trigger on a view in PL/pgSQL.

Example 44.5. A PL/pgSQL View Trigger Function for Auditing

This example uses a trigger on the view to make it updatable, and ensure that any insert, update or delete of a row in the view is recorded (i.e., audited) in the `emp_audit` table. The current time and user name are recorded, together with the type of operation performed, and the view displays the last modified time of each row.

```
CREATE TABLE emp (
    empname      text PRIMARY KEY,
```

```

        salary            integer
    );

CREATE TABLE emp_audit(
    operation             char(1)    NOT NULL,
    userid                text       NOT NULL,
    empname                text       NOT NULL,
    salary                integer,
    stamp                 timestamp NOT NULL
);

CREATE VIEW emp_view AS
    SELECT e.empname,
           e.salary,
           max(ea.stamp) AS last_updated
    FROM emp e
    LEFT JOIN emp_audit ea ON ea.empname = e.empname
    GROUP BY 1, 2;

CREATE OR REPLACE FUNCTION update_emp_view() RETURNS TRIGGER AS $$
BEGIN
    --
    -- Perform the required operation on emp, and create a row in emp_audit
    -- to reflect the change made to emp.
    --
    IF (TG_OP = 'DELETE') THEN
        DELETE FROM emp WHERE empname = OLD.empname;
        IF NOT FOUND THEN RETURN NULL; END IF;

        OLD.last_updated = now();
        INSERT INTO emp_audit VALUES('D', current_user, OLD.*);
        RETURN OLD;
    ELSEIF (TG_OP = 'UPDATE') THEN
        UPDATE emp SET salary = NEW.salary WHERE empname = OLD.empname;
        IF NOT FOUND THEN RETURN NULL; END IF;

        NEW.last_updated = now();
        INSERT INTO emp_audit VALUES('U', current_user, NEW.*);
        RETURN NEW;
    ELSEIF (TG_OP = 'INSERT') THEN
        INSERT INTO emp VALUES(NEW.empname, NEW.salary);

        NEW.last_updated = now();
        INSERT INTO emp_audit VALUES('I', current_user, NEW.*);
        RETURN NEW;
    END IF;
END;
$$ LANGUAGE plpgsql;

CREATE TRIGGER emp_audit
INSTEAD OF INSERT OR UPDATE OR DELETE ON emp_view
FOR EACH ROW EXECUTE FUNCTION update_emp_view();

```

One use of triggers is to maintain a summary table of another table. The resulting summary can be used in place of the original table for certain queries — often with vastly reduced run times. This technique is commonly used in Data Warehousing, where the tables of measured or observed data (called fact tables) might be extremely large. [Example 44.6](#) shows an example of a trigger function in PL/pgSQL that maintains a summary table for a fact table in a data warehouse.

Example 44.6. A PL/pgSQL Trigger Function for Maintaining a Summary Table

The schema detailed here is partly based on the *Grocery Store* example from *The Data Warehouse Toolkit* by Ralph Kimball.

```
--
-- Main tables - time dimension and sales fact.
--
CREATE TABLE time_dimension (
    time_key          integer NOT NULL,
    day_of_week       integer NOT NULL,
    day_of_month       integer NOT NULL,
    month             integer NOT NULL,
    quarter           integer NOT NULL,
    year              integer NOT NULL
);
CREATE UNIQUE INDEX time_dimension_key ON time_dimension(time_key);

CREATE TABLE sales_fact (
    time_key          integer NOT NULL,
    product_key       integer NOT NULL,
    store_key         integer NOT NULL,
    amount_sold       numeric(12,2) NOT NULL,
    units_sold        integer NOT NULL,
    amount_cost       numeric(12,2) NOT NULL
);
CREATE INDEX sales_fact_time ON sales_fact(time_key);

--
-- Summary table - sales by time.
--
CREATE TABLE sales_summary_bytime (
    time_key          integer NOT NULL,
    amount_sold       numeric(15,2) NOT NULL,
    units_sold        numeric(12) NOT NULL,
    amount_cost       numeric(15,2) NOT NULL
);
CREATE UNIQUE INDEX sales_summary_bytime_key ON sales_summary_bytime(time_key);

--
-- Function and trigger to amend summarized column(s) on UPDATE, INSERT, DELETE.
--
CREATE OR REPLACE FUNCTION maint_sales_summary_bytime() RETURNS TRIGGER
AS $maint_sales_summary_bytime$
    DECLARE
        delta_time_key          integer;
        delta_amount_sold       numeric(15,2);
        delta_units_sold        numeric(12);
        delta_amount_cost       numeric(15,2);
    BEGIN

        -- Work out the increment/decrement amount(s).
        IF (TG_OP = 'DELETE') THEN

            delta_time_key = OLD.time_key;
            delta_amount_sold = -1 * OLD.amount_sold;
            delta_units_sold = -1 * OLD.units_sold;
            delta_amount_cost = -1 * OLD.amount_cost;
```

```
ELSIF (TG_OP = 'UPDATE') THEN

    -- forbid updates that change the time_key -
    -- (probably not too onerous, as DELETE + INSERT is how most
    -- changes will be made).
    IF ( OLD.time_key != NEW.time_key) THEN
        RAISE EXCEPTION 'Update of time_key : % -> % not allowed',
                        OLD.time_key, NEW.time_key;
    END IF;

    delta_time_key = OLD.time_key;
    delta_amount_sold = NEW.amount_sold - OLD.amount_sold;
    delta_units_sold = NEW.units_sold - OLD.units_sold;
    delta_amount_cost = NEW.amount_cost - OLD.amount_cost;

ELSIF (TG_OP = 'INSERT') THEN

    delta_time_key = NEW.time_key;
    delta_amount_sold = NEW.amount_sold;
    delta_units_sold = NEW.units_sold;
    delta_amount_cost = NEW.amount_cost;

END IF;

-- Insert or update the summary row with the new values.
<<insert_update>>
LOOP
    UPDATE sales_summary_bytime
        SET amount_sold = amount_sold + delta_amount_sold,
            units_sold = units_sold + delta_units_sold,
            amount_cost = amount_cost + delta_amount_cost
        WHERE time_key = delta_time_key;

    EXIT insert_update WHEN found;

BEGIN
    INSERT INTO sales_summary_bytime (
        time_key,
        amount_sold,
        units_sold,
        amount_cost)
        VALUES (
            delta_time_key,
            delta_amount_sold,
            delta_units_sold,
            delta_amount_cost
        );

    EXIT insert_update;

EXCEPTION
    WHEN UNIQUE_VIOLATION THEN
        -- do nothing
END;
END LOOP insert_update;

RETURN NULL;
```

```

END;
$maint_sales_summary_bytime$ LANGUAGE plpgsql;

CREATE TRIGGER maint_sales_summary_bytime
AFTER INSERT OR UPDATE OR DELETE ON sales_fact
    FOR EACH ROW EXECUTE FUNCTION maint_sales_summary_bytime();

INSERT INTO sales_fact VALUES (1,1,1,10,3,15);
INSERT INTO sales_fact VALUES (1,2,1,20,5,35);
INSERT INTO sales_fact VALUES (2,2,1,40,15,135);
INSERT INTO sales_fact VALUES (2,3,1,10,1,13);
SELECT * FROM sales_summary_bytime;
DELETE FROM sales_fact WHERE product_key = 1;
SELECT * FROM sales_summary_bytime;
UPDATE sales_fact SET units_sold = units_sold * 2;
SELECT * FROM sales_summary_bytime;

```

AFTER triggers can also make use of *transition tables* to inspect the entire set of rows changed by the triggering statement. The CREATE TRIGGER command assigns names to one or both transition tables, and then the function can refer to those names as though they were read-only temporary tables. [Example 44.7](#) shows an example.

Example 44.7. Auditing with Transition Tables

This example produces the same results as [Example 44.4](#), but instead of using a trigger that fires for every row, it uses a trigger that fires once per statement, after collecting the relevant information in a transition table. This can be significantly faster than the row-trigger approach when the invoking statement has modified many rows. Notice that we must make a separate trigger declaration for each kind of event, since the REFERENCING clauses must be different for each case. But this does not stop us from using a single trigger function if we choose. (In practice, it might be better to use three separate functions and avoid the run-time tests on TG_OP.)

```

CREATE TABLE emp (
    empname      text NOT NULL,
    salary       integer
);

CREATE TABLE emp_audit (
    operation     char(1) NOT NULL,
    stamp        timestamp NOT NULL,
    userid       text NOT NULL,
    empname      text NOT NULL,
    salary       integer
);

CREATE OR REPLACE FUNCTION process_emp_audit() RETURNS TRIGGER AS $emp_audit$
BEGIN
    --
    -- Create rows in emp_audit to reflect the operations performed on emp,
    -- making use of the special variable TG_OP to work out the operation.
    --
    IF (TG_OP = 'DELETE') THEN
        INSERT INTO emp_audit
            SELECT 'D', now(), current_user, o.* FROM old_table o;
    ELSIF (TG_OP = 'UPDATE') THEN
        INSERT INTO emp_audit
            SELECT 'U', now(), current_user, n.* FROM new_table n;
    ELSIF (TG_OP = 'INSERT') THEN

```

```

        INSERT INTO emp_audit
            SELECT 'I', now(), current_user, n.* FROM new_table n;
    END IF;
    RETURN NULL; -- result is ignored since this is an AFTER trigger
END;
$emp_audit$ LANGUAGE plpgsql;

CREATE TRIGGER emp_audit_ins
    AFTER INSERT ON emp
    REFERENCING NEW TABLE AS new_table
    FOR EACH STATEMENT EXECUTE FUNCTION process_emp_audit();
CREATE TRIGGER emp_audit_upd
    AFTER UPDATE ON emp
    REFERENCING OLD TABLE AS old_table NEW TABLE AS new_table
    FOR EACH STATEMENT EXECUTE FUNCTION process_emp_audit();
CREATE TRIGGER emp_audit_del
    AFTER DELETE ON emp
    REFERENCING OLD TABLE AS old_table
    FOR EACH STATEMENT EXECUTE FUNCTION process_emp_audit();

```

44.10.2. Triggers on Events

PL/pgSQL can be used to define [event triggers](#). Postgres Pro requires that a function that is to be called as an event trigger must be declared as a function with no arguments and a return type of `event_trigger`.

When a PL/pgSQL function is called as an event trigger, several special variables are created automatically in the top-level block. They are:

`TG_EVENT` text

event the trigger is fired for.

`TG_TAG` text

command tag for which the trigger is fired.

[Example 44.8](#) shows an example of an event trigger function in PL/pgSQL.

Example 44.8. A PL/pgSQL Event Trigger Function

This example trigger simply raises a `NOTICE` message each time a supported command is executed.

```

CREATE OR REPLACE FUNCTION snitch() RETURNS event_trigger AS $$
BEGIN
    RAISE NOTICE 'snitch: % %', tg_event, tg_tag;
END;
$$ LANGUAGE plpgsql;

CREATE EVENT TRIGGER snitch ON ddl_command_start EXECUTE FUNCTION snitch();

```

44.11. Packages

PL/pgSQL can be used to define functions for working with *packages*. A package is essentially a schema that helps to organize the named objects with a related purpose so it can be created using [CREATE SCHEMA](#) or a special [CREATE PACKAGE](#) command. A package should contain only functions, procedures and composite types.

When you create a package, you must also create the initialization function. The initialization function is a PL/pgSQL function named `__init__` that has no arguments and returns `void`. Note that the initialization function must be defined before any other package function, and all the variables declared in the

initialization function are global so they can be used by functions of other packages using dot notation. For example, the `bar` variable declared in the `__init__` function of the `foo` package can be referenced as `foo.bar`.

Postgres Pro provides function creation modifiers for working with packages. The `#package` modifier defines that it is a package function and it can use variables of this package directly. The `#import` modifier defines that the function can use variables of other packages using the dot notation as described above. The `#package` modifier is omitted when creating functions inside a package using the `CREATE PACKAGE` command. The `#private` modifier indicates that the function is private, meaning that it cannot be referenced from outside the package, but it is necessary for the internal processes of the package. The modifier `#export` indicates that the package variable is public, meaning it can be referenced from outside the package. When a function with the `#package` modifier is called, the containing package is initialized if it has not been initialized already in the current session. You can create a function with the `#package` modifier only in the schema containing the initialization function. If the `#import` modifier is specified, the function will have access to the variables of the packages specified in a comma-separated list. These packages are initialized automatically if they have not been initialized already in the current session. In packages created with `CREATE PACKAGE`, `#import` specified for the initialization function affects all package functions, while in `CREATE SCHEMA` it must be defined separately for `__init__` and other package functions. The `#import` modifier can also be used in anonymous code blocks.

In the below example of using the `#import` modifier, the `showValues` procedure calls `p` of the `http` package and `set_action` of the `dbms_application_info` package.

```
CREATE OR REPLACE PROCEDURE showValues(p_Str varchar) AS $$
#import http, dbms_application_info
BEGIN
    CALL dbms_application_info.set_action('Show hello');

    CALL http.p('<p>' || p_Str || '</p>');
END;
$$LANGUAGE plpgsql;
```

PL/pgSQL also provides a built-in function `plpgsql_reset_packages()` that deinstantiates all the packages in this session to bring the session to the original state.

```
SELECT plpgsql_reset_packages();
```

44.12. PL/pgSQL under the Hood

This section discusses some implementation details that are frequently important for PL/pgSQL users to know.

44.12.1. Variable Substitution

SQL statements and expressions within a PL/pgSQL function can refer to variables and parameters of the function. Behind the scenes, PL/pgSQL substitutes query parameters for such references. Query parameters will only be substituted in places where they are syntactically permissible. As an extreme case, consider this example of poor programming style:

```
INSERT INTO foo (foo) VALUES (foo(foo));
```

The first occurrence of `foo` must syntactically be a table name, so it will not be substituted, even if the function has a variable named `foo`. The second occurrence must be the name of a column of that table, so it will not be substituted either. Likewise the third occurrence must be a function name, so it also will not be substituted for. Only the last occurrence is a candidate to be a reference to a variable of the PL/pgSQL function.

Another way to understand this is that variable substitution can only insert data values into an SQL command; it cannot dynamically change which database objects are referenced by the command. (If you want to do that, you must build a command string dynamically, as explained in [Section 44.5.4.](#))

Since the names of variables are syntactically no different from the names of table columns, there can be ambiguity in statements that also refer to tables: is a given name meant to refer to a table column, or a variable? Let's change the previous example to

```
INSERT INTO dest (col) SELECT foo + bar FROM src;
```

Here, `dest` and `src` must be table names, and `col` must be a column of `dest`, but `foo` and `bar` might reasonably be either variables of the function or columns of `src`.

By default, PL/pgSQL will report an error if a name in an SQL statement could refer to either a variable or a table column. You can fix such a problem by renaming the variable or column, or by qualifying the ambiguous reference, or by telling PL/pgSQL which interpretation to prefer.

The simplest solution is to rename the variable or column. A common coding rule is to use a different naming convention for PL/pgSQL variables than you use for column names. For example, if you consistently name function variables `v_something` while none of your column names start with `v_`, no conflicts will occur.

Alternatively you can qualify ambiguous references to make them clear. In the above example, `src.foo` would be an unambiguous reference to the table column. To create an unambiguous reference to a variable, declare it in a labeled block and use the block's label (see [Section 44.2](#)). For example,

```
<<block>>
DECLARE
    foo int;
BEGIN
    foo := ...;
    INSERT INTO dest (col) SELECT block.foo + bar FROM src;
```

Here `block.foo` means the variable even if there is a column `foo` in `src`. Function parameters, as well as special variables such as `FOUND`, can be qualified by the function's name, because they are implicitly declared in an outer block labeled with the function's name.

Sometimes it is impractical to fix all the ambiguous references in a large body of PL/pgSQL code. In such cases you can specify that PL/pgSQL should resolve ambiguous references as the variable (which is compatible with PL/pgSQL's behavior before PostgreSQL 9.0), or as the table column (which is compatible with some other systems such as Oracle).

To change this behavior on a system-wide basis, set the configuration parameter `plpgsql.variable_conflict` to one of `error`, `use_variable`, or `use_column` (where `error` is the factory default). This parameter affects subsequent compilations of statements in PL/pgSQL functions, but not statements already compiled in the current session. Because changing this setting can cause unexpected changes in the behavior of PL/pgSQL functions, it can only be changed by a superuser.

You can also set the behavior on a function-by-function basis, by inserting one of these special commands at the start of the function text:

```
#variable_conflict error
#variable_conflict use_variable
#variable_conflict use_column
```

These commands affect only the function they are written in, and override the setting of `plpgsql.variable_conflict`. An example is

```
CREATE FUNCTION stamp_user(id int, comment text) RETURNS void AS $$
    #variable_conflict use_variable
    DECLARE
        curtime timestamp := now();
    BEGIN
        UPDATE users SET last_modified = curtime, comment = comment
            WHERE users.id = id;
    END;
$$ LANGUAGE plpgsql;
```

In the `UPDATE` command, `curtime`, `comment`, and `id` will refer to the function's variable and parameters whether or not `users` has columns of those names. Notice that we had to qualify the reference to `users.id` in the `WHERE` clause to make it refer to the table column. But we did not have to qualify the reference to `comment` as a target in the `UPDATE` list, because syntactically that must be a column of `users`. We could write the same function without depending on the `variable_conflict` setting in this way:

```
CREATE FUNCTION stamp_user(id int, comment text) RETURNS void AS $$
    <<fn>>
    DECLARE
        curtime timestamp := now();
    BEGIN
        UPDATE users SET last_modified = fn.curtime, comment = stamp_user.comment
            WHERE users.id = stamp_user.id;
    END;
$$ LANGUAGE plpgsql;
```

Variable substitution does not happen in a command string given to `EXECUTE` or one of its variants. If you need to insert a varying value into such a command, do so as part of constructing the string value, or use `USING`, as illustrated in [Section 44.5.4](#).

Variable substitution currently works only in `SELECT`, `INSERT`, `UPDATE`, `DELETE`, and commands containing one of these (such as `EXPLAIN` and `CREATE TABLE ... AS SELECT`), because the main SQL engine allows query parameters only in these commands. To use a non-constant name or value in other statement types (generically called utility statements), you must construct the utility statement as a string and `EXECUTE` it.

44.12.2. Plan Caching

The PL/pgSQL interpreter parses the function's source text and produces an internal binary instruction tree the first time the function is called (within each session). The instruction tree fully translates the PL/pgSQL statement structure, but individual SQL expressions and SQL commands used in the function are not translated immediately.

As each expression and SQL command is first executed in the function, the PL/pgSQL interpreter parses and analyzes the command to create a prepared statement, using the SPI manager's `SPI_prepare` function. Subsequent visits to that expression or command reuse the prepared statement. Thus, a function with conditional code paths that are seldom visited will never incur the overhead of analyzing those commands that are never executed within the current session. A disadvantage is that errors in a specific expression or command cannot be detected until that part of the function is reached in execution. (Trivial syntax errors will be detected during the initial parsing pass, but anything deeper will not be detected until execution.)

PL/pgSQL (or more precisely, the SPI manager) can furthermore attempt to cache the execution plan associated with any particular prepared statement. If a cached plan is not used, then a fresh execution plan is generated on each visit to the statement, and the current parameter values (that is, PL/pgSQL variable values) can be used to optimize the selected plan. If the statement has no parameters, or is executed many times, the SPI manager will consider creating a *generic* plan that is not dependent on specific parameter values, and caching that for re-use. Typically this will happen only if the execution plan is not very sensitive to the values of the PL/pgSQL variables referenced in it. If it is, generating a plan each time is a net win. See [PREPARE](#) for more information about the behavior of prepared statements.

Because PL/pgSQL saves prepared statements and sometimes execution plans in this way, SQL commands that appear directly in a PL/pgSQL function must refer to the same tables and columns on every execution; that is, you cannot use a parameter as the name of a table or column in an SQL command. To get around this restriction, you can construct dynamic commands using the PL/pgSQL `EXECUTE` statement — at the price of performing new parse analysis and constructing a new execution plan on every execution.

The mutable nature of record variables presents another problem in this connection. When fields of a record variable are used in expressions or statements, the data types of the fields must not change from

one call of the function to the next, since each expression will be analyzed using the data type that is present when the expression is first reached. `EXECUTE` can be used to get around this problem when necessary.

If the same function is used as a trigger for more than one table, PL/pgSQL prepares and caches statements independently for each such table — that is, there is a cache for each trigger function and table combination, not just for each function. This alleviates some of the problems with varying data types; for instance, a trigger function will be able to work successfully with a column named `key` even if it happens to have different types in different tables.

Likewise, functions having polymorphic argument types have a separate statement cache for each combination of actual argument types they have been invoked for, so that data type differences do not cause unexpected failures.

Statement caching can sometimes have surprising effects on the interpretation of time-sensitive values. For example there is a difference between what these two functions do:

```
CREATE FUNCTION logfunc1(logtxt text) RETURNS void AS $$
BEGIN
    INSERT INTO logtable VALUES (logtxt, 'now');
END;
$$ LANGUAGE plpgsql;
```

and:

```
CREATE FUNCTION logfunc2(logtxt text) RETURNS void AS $$
DECLARE
    curtime timestamp;
BEGIN
    curtime := 'now';
    INSERT INTO logtable VALUES (logtxt, curtime);
END;
$$ LANGUAGE plpgsql;
```

In the case of `logfunc1`, the Postgres Pro main parser knows when analyzing the `INSERT` that the string `'now'` should be interpreted as `timestamp`, because the target column of `logtable` is of that type. Thus, `'now'` will be converted to a `timestamp` constant when the `INSERT` is analyzed, and then used in all invocations of `logfunc1` during the lifetime of the session. Needless to say, this isn't what the programmer wanted. A better idea is to use the `now()` or `current_timestamp` function.

In the case of `logfunc2`, the Postgres Pro main parser does not know what type `'now'` should become and therefore it returns a data value of type `text` containing the string `now`. During the ensuing assignment to the local variable `curtime`, the PL/pgSQL interpreter casts this string to the `timestamp` type by calling the `textout` and `timestamp_in` functions for the conversion. So, the computed time stamp is updated on each execution as the programmer expects. Even though this happens to work as expected, it's not terribly efficient, so use of the `now()` function would still be a better idea.

44.13. Tips for Developing in PL/pgSQL

One good way to develop in PL/pgSQL is to use the text editor of your choice to create your functions, and in another window, use `psql` to load and test those functions. If you are doing it this way, it is a good idea to write the function using `CREATE OR REPLACE FUNCTION`. That way you can just reload the file to update the function definition. For example:

```
CREATE OR REPLACE FUNCTION testfunc(integer) RETURNS integer AS $$
    ....
$$ LANGUAGE plpgsql;
```

While running `psql`, you can load or reload such a function definition file with:

```
\i filename.sql
```


and then immediately issue SQL commands to test the function.

Another good way to develop in PL/pgSQL is with a GUI database access tool that facilitates development in a procedural language. One example of such a tool is pgAdmin, although others exist. These tools often provide convenient features such as escaping single quotes and making it easier to recreate and debug functions.

44.13.1. Handling of Quotation Marks

The code of a PL/pgSQL function is specified in `CREATE FUNCTION` as a string literal. If you write the string literal in the ordinary way with surrounding single quotes, then any single quotes inside the function body must be doubled; likewise any backslashes must be doubled (assuming escape string syntax is used). Doubling quotes is at best tedious, and in more complicated cases the code can become downright incomprehensible, because you can easily find yourself needing half a dozen or more adjacent quote marks. It's recommended that you instead write the function body as a “dollar-quoted” string literal (see [Section 4.1.2.4](#)). In the dollar-quoting approach, you never double any quote marks, but instead take care to choose a different dollar-quoting delimiter for each level of nesting you need. For example, you might write the `CREATE FUNCTION` command as:

```
CREATE OR REPLACE FUNCTION testfunc(integer) RETURNS integer AS $PROC$
    ....
$PROC$ LANGUAGE plpgsql;
```

Within this, you might use quote marks for simple literal strings in SQL commands and `$$` to delimit fragments of SQL commands that you are assembling as strings. If you need to quote text that includes `$$`, you could use `Q`, and so on.

The following chart shows what you have to do when writing quote marks without dollar quoting. It might be useful when translating pre-dollar quoting code into something more comprehensible.

1 quotation mark

To begin and end the function body, for example:

```
CREATE FUNCTION foo() RETURNS integer AS '
    ....
' LANGUAGE plpgsql;
```

Anywhere within a single-quoted function body, quote marks *must* appear in pairs.

2 quotation marks

For string literals inside the function body, for example:

```
a_output := 'Blah';
SELECT * FROM users WHERE f_name='foobar';
```

In the dollar-quoting approach, you'd just write:

```
a_output := 'Blah';
SELECT * FROM users WHERE f_name='foobar';
```

which is exactly what the PL/pgSQL parser would see in either case.

4 quotation marks

When you need a single quotation mark in a string constant inside the function body, for example:

```
a_output := a_output || ' AND name LIKE '''foobar''' AND xyz'
```

The value actually appended to `a_output` would be: `AND name LIKE 'foobar' AND xyz`.

In the dollar-quoting approach, you'd write:

```
a_output := a_output || $$ AND name LIKE 'foobar' AND xyz$$
```

being careful that any dollar-quote delimiters around this are not just \$\$.

6 quotation marks

When a single quotation mark in a string inside the function body is adjacent to the end of that string constant, for example:

```
a_output := a_output || ' AND name LIKE '''foobar'''''
```

The value appended to `a_output` would then be: `AND name LIKE 'foobar'`.

In the dollar-quoting approach, this becomes:

```
a_output := a_output || $$ AND name LIKE 'foobar'$$
```

10 quotation marks

When you want two single quotation marks in a string constant (which accounts for 8 quotation marks) and this is adjacent to the end of that string constant (2 more). You will probably only need that if you are writing a function that generates other functions, as in [Example 44.10](#). For example:

```

a_output := a_output || ' if v_' ||
referrer_keys.kind || ' like ' ||
|| referrer_keys.key_string || ' ' ||
then return ' ' || referrer_keys.referrer_type
|| ' ' ; end if ;

```

The value of `a_output` would then be:

```
if v_... like '...' then return '...'; end if;
```

In the dollar-quoting approach, this becomes:

```
a_output := a_output || $$ if v_$$ || referrer_keys.kind || $$ like '$$
|| referrer_keys.key_string || $$'
then return '$$ || referrer_keys.referrer_type
|| $$'; end if;$$;
```

where we assume we only need to put single quote marks into `a_output`, because it will be re-quoted before use.

44.13.2. Additional Compile-Time and Run-Time Checks

To aid the user in finding instances of simple but common problems before they cause harm, PL/pgSQL provides additional *checks*. When enabled, depending on the configuration, they can be used to emit either a `WARNING` or an `ERROR` during the compilation of a function. A function which has received a `WARNING` can be executed without producing further messages, so you are advised to test in a separate development environment.

Setting `plpgsql.extra_warnings`, or `plpgsql.extra_errors`, as appropriate, to "all" is encouraged in development and/or testing environments.

These additional checks are enabled through the configuration variables `plpgsql.extra_warnings` for warnings and `plpgsql.extra_errors` for errors. Both can be set either to a comma-separated list of checks, "none" or "all". The default is "none". Currently the list of available checks includes:

shadowed variables

Checks if a declaration shadows a previously defined variable.

```
strict_multi_assignment
```

Some PL/pgSQL commands allow assigning values to more than one variable at a time, such as `SELECT INTO`. Typically, the number of target variables and the number of source variables should match, though PL/pgSQL will use `NULL` for missing values and extra variables are ignored. Enabling

this check will cause PL/pgSQL to throw a `WARNING` or `ERROR` whenever the number of target variables and the number of source variables are different.

`too_many_rows`

Enabling this check will cause PL/pgSQL to check if a given query returns more than one row when an `INTO` clause is used. As an `INTO` statement will only ever use one row, having a query return multiple rows is generally either inefficient and/or nondeterministic and therefore is likely an error.

The following example shows the effect of `plpgsql.extra_warnings` set to `shadowed_variables`:

```
SET plpgsql.extra_warnings TO 'shadowed_variables';

CREATE FUNCTION foo(f1 int) RETURNS int AS $$
DECLARE
f1 int;
BEGIN
RETURN f1;
END;
$$ LANGUAGE plpgsql;
WARNING: variable "f1" shadows a previously defined variable
LINE 3: f1 int;
      ^
CREATE FUNCTION
```

The below example shows the effects of setting `plpgsql.extra_warnings` to `strict_multi_assignment`:

```
SET plpgsql.extra_warnings TO 'strict_multi_assignment';

CREATE OR REPLACE FUNCTION public.foo()
RETURNS void
LANGUAGE plpgsql
AS $$
DECLARE
    x int;
    y int;
BEGIN
    SELECT 1 INTO x, y;
    SELECT 1, 2 INTO x, y;
    SELECT 1, 2, 3 INTO x, y;
END;
$$;

SELECT foo();
WARNING: number of source and target fields in assignment does not match
DETAIL: strict_multi_assignment check of extra_warnings is active.
HINT: Make sure the query returns the exact list of columns.
WARNING: number of source and target fields in assignment does not match
DETAIL: strict_multi_assignment check of extra_warnings is active.
HINT: Make sure the query returns the exact list of columns.

foo
----

(1 row)
```

44.14. Porting from Oracle PL/SQL

This section explains differences between Postgres Pro's PL/pgSQL language and Oracle's PL/SQL language, to help developers who port applications from Oracle® to Postgres Pro.

PL/pgSQL is similar to PL/SQL in many aspects. It is a block-structured, imperative language, and all variables have to be declared. Assignments, loops, and conditionals are similar. The main differences you should keep in mind when porting from PL/SQL to PL/pgSQL are:

- If a name used in an SQL command could be either a column name of a table used in the command or a reference to a variable of the function, PL/SQL treats it as a column name. By default, PL/pgSQL will throw an error complaining that the name is ambiguous. You can specify `plpgsql.variable_conflict = use_column` to change this behavior to match PL/SQL, as explained in [Section 44.12.1](#). It's often best to avoid such ambiguities in the first place, but if you have to port a large amount of code that depends on this behavior, setting `variable_conflict` may be the best solution.
- In Postgres Pro the function body must be written as a string literal. Therefore you need to use dollar quoting or escape single quotes in the function body. (See [Section 44.13.1](#).)
- Data type names often need translation. For example, in Oracle string values are commonly declared as being of type `varchar2`, which is a non-SQL-standard type. In Postgres Pro, use type `varchar` or `text` instead. Similarly, replace type `number` with `numeric`, or use some other numeric data type if there's a more appropriate one.
- Integer `FOR` loops with `REVERSE` work differently: PL/SQL counts down from the second number to the first, while PL/pgSQL counts down from the first number to the second, requiring the loop bounds to be swapped when porting. This incompatibility is unfortunate but is unlikely to be changed. (See [Section 44.6.5.5](#).)
- `FOR` loops over queries (other than cursors) also work differently: the target variable(s) must have been declared, whereas PL/SQL always declares them implicitly. An advantage of this is that the variable values are still accessible after the loop exits.
- There are various notational differences for the use of cursor variables.

44.14.1. Porting Examples

[Example 44.9](#) shows how to port a simple function from PL/SQL to PL/pgSQL.

Example 44.9. Porting a Simple Function from PL/SQL to PL/pgSQL

Here is an Oracle PL/SQL function:

```
CREATE OR REPLACE FUNCTION cs_fmt_browser_version(v_name varchar2,
                                                    v_version varchar2)
RETURN varchar2 IS
BEGIN
    IF v_version IS NULL THEN
        RETURN v_name;
    END IF;
    RETURN v_name || '/' || v_version;
END;
/
show errors;
```

Let's go through this function and see the differences compared to PL/pgSQL:

- The type name `varchar2` has to be changed to `varchar` or `text`. In the examples in this section, we'll use `varchar`, but `text` is often a better choice if you do not need specific string length limits.
- The `RETURN` key word in the function prototype (not the function body) becomes `RETURNS` in Postgres Pro. Also, `IS` becomes `AS`, and you need to add a `LANGUAGE` clause because PL/pgSQL is not the only possible function language.
- In Postgres Pro, the function body is considered to be a string literal, so you need to use quote marks or dollar quotes around it. This substitutes for the terminating `/` in the Oracle approach.

- The `show errors` command does not exist in Postgres Pro, and is not needed since errors are reported automatically.

This is how this function would look when ported to Postgres Pro:

```
CREATE OR REPLACE FUNCTION cs_fmt_browser_version(v_name varchar,
                                                  v_version varchar)
RETURNS varchar AS $$
BEGIN
    IF v_version IS NULL THEN
        RETURN v_name;
    END IF;
    RETURN v_name || '/' || v_version;
END;
$$ LANGUAGE plpgsql;
```

[Example 44.10](#) shows how to port a function that creates another function and how to handle the ensuing quoting problems.

Example 44.10. Porting a Function that Creates Another Function from PL/SQL to PL/pgSQL

The following procedure grabs rows from a `SELECT` statement and builds a large function with the results in `IF` statements, for the sake of efficiency.

This is the Oracle version:

```
CREATE OR REPLACE PROCEDURE cs_update_referrer_type_proc IS
    CURSOR referrer_keys IS
        SELECT * FROM cs_referrer_keys
        ORDER BY try_order;
    func_cmd VARCHAR(4000);
BEGIN
    func_cmd := 'CREATE OR REPLACE FUNCTION cs_find_referrer_type(v_host IN VARCHAR2,
        v_domain IN VARCHAR2, v_url IN VARCHAR2) RETURN VARCHAR2 IS BEGIN';

    FOR referrer_key IN referrer_keys LOOP
        func_cmd := func_cmd ||
            ' IF v_' || referrer_key.kind
            || ' LIKE ''' || referrer_key.key_string
            || ''' THEN RETURN ''' || referrer_key.referrer_type
            || '''; END IF;';
    END LOOP;

    func_cmd := func_cmd || ' RETURN NULL; END;';

    EXECUTE IMMEDIATE func_cmd;
END;
/
show errors;
```

Here is how this function would end up in Postgres Pro:

```
CREATE OR REPLACE PROCEDURE cs_update_referrer_type_proc() AS $func$
DECLARE
    referrer_keys CURSOR IS
        SELECT * FROM cs_referrer_keys
        ORDER BY try_order;
    func_body text;
    func_cmd text;
BEGIN
    func_body := 'BEGIN';
```

```

FOR referrer_key IN referrer_keys LOOP
    func_body := func_body ||
        ' IF v_' || referrer_key.kind
        || ' LIKE ' || quote_literal(referrer_key.key_string)
        || ' THEN RETURN ' || quote_literal(referrer_key.referrer_type)
        || '; END IF; ' ;
END LOOP;

func_body := func_body || ' RETURN NULL; END;';

func_cmd :=
    'CREATE OR REPLACE FUNCTION cs_find_referrer_type(v_host varchar,
                                                    v_domain varchar,
                                                    v_url varchar)
      RETURNS varchar AS '
    || quote_literal(func_body)
    || ' LANGUAGE plpgsql;' ;

EXECUTE func_cmd;
END;
$func$ LANGUAGE plpgsql;

```

Notice how the body of the function is built separately and passed through `quote_literal` to double any quote marks in it. This technique is needed because we cannot safely use dollar quoting for defining the new function: we do not know for sure what strings will be interpolated from the `referrer_key.key_string` field. (We are assuming here that `referrer_key.kind` can be trusted to always be host, domain, or url, but `referrer_key.key_string` might be anything, in particular it might contain dollar signs.) This function is actually an improvement on the Oracle original, because it will not generate broken code when `referrer_key.key_string` or `referrer_key.referrer_type` contain quote marks.

[Example 44.11](#) shows how to port a function with `OUT` parameters and string manipulation. Postgres Pro does not have a built-in `instr` function, but you can create one using a combination of other functions. In [Section 44.14.3](#) there is a PL/pgSQL implementation of `instr` that you can use to make your porting easier.

Example 44.11. Porting a Procedure With String Manipulation and `OUT` Parameters from PL/SQL to PL/pgSQL

The following Oracle PL/SQL procedure is used to parse a URL and return several elements (host, path, and query).

This is the Oracle version:

```

CREATE OR REPLACE PROCEDURE cs_parse_url(
    v_url IN VARCHAR2,
    v_host OUT VARCHAR2,  -- This will be passed back
    v_path OUT VARCHAR2,  -- This one too
    v_query OUT VARCHAR2) -- And this one
IS
    a_pos1 INTEGER;
    a_pos2 INTEGER;
BEGIN
    v_host := NULL;
    v_path := NULL;
    v_query := NULL;
    a_pos1 := instr(v_url, '//');

    IF a_pos1 = 0 THEN
        RETURN;
    END IF;

```

```

END IF;
a_pos2 := instr(v_url, '/', a_pos1 + 2);
IF a_pos2 = 0 THEN
    v_host := substr(v_url, a_pos1 + 2);
    v_path := '/';
    RETURN;
END IF;

v_host := substr(v_url, a_pos1 + 2, a_pos2 - a_pos1 - 2);
a_pos1 := instr(v_url, '?', a_pos2 + 1);

IF a_pos1 = 0 THEN
    v_path := substr(v_url, a_pos2);
    RETURN;
END IF;

v_path := substr(v_url, a_pos2, a_pos1 - a_pos2);
v_query := substr(v_url, a_pos1 + 1);
END;
/
show errors;

```

Here is a possible translation into PL/pgSQL:

```

CREATE OR REPLACE FUNCTION cs_parse_url(
    v_url IN VARCHAR,
    v_host OUT VARCHAR, -- This will be passed back
    v_path OUT VARCHAR, -- This one too
    v_query OUT VARCHAR) -- And this one
AS $$
DECLARE
    a_pos1 INTEGER;
    a_pos2 INTEGER;
BEGIN
    v_host := NULL;
    v_path := NULL;
    v_query := NULL;
    a_pos1 := instr(v_url, '//');

    IF a_pos1 = 0 THEN
        RETURN;
    END IF;
    a_pos2 := instr(v_url, '/', a_pos1 + 2);
    IF a_pos2 = 0 THEN
        v_host := substr(v_url, a_pos1 + 2);
        v_path := '/';
        RETURN;
    END IF;

    v_host := substr(v_url, a_pos1 + 2, a_pos2 - a_pos1 - 2);
    a_pos1 := instr(v_url, '?', a_pos2 + 1);

    IF a_pos1 = 0 THEN
        v_path := substr(v_url, a_pos2);
        RETURN;
    END IF;

    v_path := substr(v_url, a_pos2, a_pos1 - a_pos2);
    v_query := substr(v_url, a_pos1 + 1);

```

```
END;
$$ LANGUAGE plpgsql;
```

This function could be used like this:

```
SELECT * FROM cs_parse_url('http://foobar.com/query.cgi?baz');
```

[Example 44.12](#) shows how to port a procedure that uses numerous features that are specific to Oracle.

Example 44.12. Porting a Procedure from PL/SQL to PL/pgSQL

The Oracle version:

```
CREATE OR REPLACE PROCEDURE cs_create_job(v_job_id IN INTEGER) IS
    a_running_job_count INTEGER;
BEGIN
    LOCK TABLE cs_jobs IN EXCLUSIVE MODE;

    SELECT count(*) INTO a_running_job_count FROM cs_jobs WHERE end_stamp IS NULL;

    IF a_running_job_count > 0 THEN
        COMMIT; -- free lock
        raise_application_error(-20000,
            'Unable to create a new job: a job is currently running.');
```

END IF;

```
DELETE FROM cs_active_job;
INSERT INTO cs_active_job(job_id) VALUES (v_job_id);

BEGIN
    INSERT INTO cs_jobs (job_id, start_stamp) VALUES (v_job_id, now());
EXCEPTION
    WHEN dup_val_on_index THEN NULL; -- don't worry if it already exists
END;
COMMIT;
END;
/
show errors
```

This is how we could port this procedure to PL/pgSQL:

```
CREATE OR REPLACE PROCEDURE cs_create_job(v_job_id integer) AS $$
DECLARE
    a_running_job_count integer;
BEGIN
    LOCK TABLE cs_jobs IN EXCLUSIVE MODE;

    SELECT count(*) INTO a_running_job_count FROM cs_jobs WHERE end_stamp IS NULL;

    IF a_running_job_count > 0 THEN
        COMMIT; -- free lock
        RAISE EXCEPTION 'Unable to create a new job: a job is currently running'; -- 1
    END IF;

    DELETE FROM cs_active_job;
    INSERT INTO cs_active_job(job_id) VALUES (v_job_id);

    BEGIN
        INSERT INTO cs_jobs (job_id, start_stamp) VALUES (v_job_id, now());
    EXCEPTION
        WHEN unique_violation THEN -- 2
```



```

        -- don't worry if it already exists
    END;
    COMMIT;
END;
$$ LANGUAGE plpgsql;

```

- ❶ The syntax of `RAISE` is considerably different from Oracle's statement, although the basic case `RAISE exception_name` works similarly.
- ❷ The exception names supported by PL/pgSQL are different from Oracle's. The set of built-in exception names is much larger (see [Appendix A](#)). There is not currently a way to declare user-defined exception names, although you can throw user-chosen `SQLSTATE` values instead.

44.14.2. Other Things to Watch For

This section explains a few other things to watch for when porting Oracle PL/SQL functions to Postgres Pro.

44.14.2.1. Implicit Rollback after Exceptions

In PL/pgSQL, when an exception is caught by an `EXCEPTION` clause, all database changes since the block's `BEGIN` are automatically rolled back. That is, the behavior is equivalent to what you'd get in Oracle with:

```

BEGIN
    SAVEPOINT s1;
    ... code here ...
EXCEPTION
    WHEN ... THEN
        ROLLBACK TO s1;
        ... code here ...
    WHEN ... THEN
        ROLLBACK TO s1;
        ... code here ...
END;

```

If you are translating an Oracle procedure that uses `SAVEPOINT` and `ROLLBACK TO` in this style, your task is easy: just omit the `SAVEPOINT` and `ROLLBACK TO`. If you have a procedure that uses `SAVEPOINT` and `ROLLBACK TO` in a different way then some actual thought will be required.

44.14.2.2. EXECUTE

The PL/pgSQL version of `EXECUTE` works similarly to the PL/SQL version, but you have to remember to use `quote_literal` and `quote_ident` as described in [Section 44.5.4](#). Constructs of the type `EXECUTE 'SELECT * FROM $1';` will not work reliably unless you use these functions.

44.14.2.3. Optimizing PL/pgSQL Functions

Postgres Pro gives you two function creation modifiers to optimize execution: “volatility” (whether the function always returns the same result when given the same arguments) and “strictness” (whether the function returns null if any argument is null). Consult the [CREATE FUNCTION](#) reference page for details.

When making use of these optimization attributes, your `CREATE FUNCTION` statement might look something like this:

```

CREATE FUNCTION foo(...) RETURNS integer AS $$
...
$$ LANGUAGE plpgsql STRICT IMMUTABLE;

```

44.14.3. Appendix

This section contains the code for a set of Oracle-compatible `instr` functions that you can use to simplify your porting efforts.

```
--
-- instr functions that mimic Oracle's counterpart
-- Syntax: instr(string1, string2 [, n [, m]])
-- where [] denotes optional parameters.
--
-- Search string1, beginning at the nth character, for the mth occurrence
-- of string2. If n is negative, search backwards, starting at the abs(n)'th
-- character from the end of string1.
-- If n is not passed, assume 1 (search starts at first character).
-- If m is not passed, assume 1 (find first occurrence).
-- Returns starting index of string2 in string1, or 0 if string2 is not found.
--
```

```
CREATE FUNCTION instr(vchar, vchar) RETURNS integer AS $$
BEGIN
    RETURN instr($1, $2, 1);
END;
$$ LANGUAGE plpgsql STRICT IMMUTABLE;
```

```
CREATE FUNCTION instr(string vchar, string_to_search_for vchar,
                      beg_index integer)
RETURNS integer AS $$
DECLARE
    pos integer NOT NULL DEFAULT 0;
    temp_str vchar;
    beg integer;
    length integer;
    ss_length integer;
BEGIN
    IF beg_index > 0 THEN
        temp_str := substring(string FROM beg_index);
        pos := position(string_to_search_for IN temp_str);

        IF pos = 0 THEN
            RETURN 0;
        ELSE
            RETURN pos + beg_index - 1;
        END IF;
    ELSIF beg_index < 0 THEN
        ss_length := char_length(string_to_search_for);
        length := char_length(string);
        beg := length + 1 + beg_index;

        WHILE beg > 0 LOOP
            temp_str := substring(string FROM beg FOR ss_length);
            IF string_to_search_for = temp_str THEN
                RETURN beg;
            END IF;

            beg := beg - 1;
        END LOOP;

        RETURN 0;
    ELSE
        RETURN 0;
    END IF;
END IF;
```

```
END;
$$ LANGUAGE plpgsql STRICT IMMUTABLE;

CREATE FUNCTION instr(string varchar, string_to_search_for varchar,
                      beg_index integer, occur_index integer)
RETURNS integer AS $$
DECLARE
    pos integer NOT NULL DEFAULT 0;
    occur_number integer NOT NULL DEFAULT 0;
    temp_str varchar;
    beg integer;
    i integer;
    length integer;
    ss_length integer;
BEGIN
    IF occur_index <= 0 THEN
        RAISE 'argument ''%'' is out of range', occur_index
        USING ERRCODE = '22003';
    END IF;

    IF beg_index > 0 THEN
        beg := beg_index - 1;
        FOR i IN 1..occur_index LOOP
            temp_str := substring(string FROM beg + 1);
            pos := position(string_to_search_for IN temp_str);
            IF pos = 0 THEN
                RETURN 0;
            END IF;
            beg := beg + pos;
        END LOOP;

        RETURN beg;
    ELSIF beg_index < 0 THEN
        ss_length := char_length(string_to_search_for);
        length := char_length(string);
        beg := length + 1 + beg_index;

        WHILE beg > 0 LOOP
            temp_str := substring(string FROM beg FOR ss_length);
            IF string_to_search_for = temp_str THEN
                occur_number := occur_number + 1;
                IF occur_number = occur_index THEN
                    RETURN beg;
                END IF;
            END IF;
            beg := beg - 1;
        END LOOP;

        RETURN 0;
    ELSE
        RETURN 0;
    END IF;
END;
$$ LANGUAGE plpgsql STRICT IMMUTABLE;
```

Chapter 45. PL/Tcl — Tcl Procedural Language

PL/Tcl is a loadable procedural language for the Postgres Pro database system that enables the [Tcl language](#) to be used to write Postgres Pro functions and procedures.

45.1. Overview

PL/Tcl offers most of the capabilities a function writer has in the C language, with a few restrictions, and with the addition of the powerful string processing libraries that are available for Tcl.

One compelling *good* restriction is that everything is executed from within the safety of the context of a Tcl interpreter. In addition to the limited command set of safe Tcl, only a few commands are available to access the database via SPI and to raise messages via `elog()`. PL/Tcl provides no way to access internals of the database server or to gain OS-level access under the permissions of the Postgres Pro server process, as a C function can do. Thus, unprivileged database users can be trusted to use this language; it does not give them unlimited authority.

The other notable implementation restriction is that Tcl functions cannot be used to create input/output functions for new data types.

Sometimes it is desirable to write Tcl functions that are not restricted to safe Tcl. For example, one might want a Tcl function that sends email. To handle these cases, there is a variant of PL/Tcl called `PL/TclU` (for untrusted Tcl). This is exactly the same language except that a full Tcl interpreter is used. *If PL/TclU is used, it must be installed as an untrusted procedural language* so that only database superusers can create functions in it. The writer of a PL/TclU function must take care that the function cannot be used to do anything unwanted, since it will be able to do anything that could be done by a user logged in as the database administrator.

The shared object code for the PL/Tcl and PL/TclU call handlers is automatically built and installed in the Postgres Pro library directory if Tcl support is specified in the configuration step of the installation procedure. To install PL/Tcl and/or PL/TclU in a particular database, use the `CREATE EXTENSION` command, for example `CREATE EXTENSION pltcl` or `CREATE EXTENSION pltclu`.

45.2. PL/Tcl Functions and Arguments

To create a function in the PL/Tcl language, use the standard [CREATE FUNCTION](#) syntax:

```
CREATE FUNCTION funcname (argument-types) RETURNS return-type AS $$
    # PL/Tcl function body
$$ LANGUAGE pltcl;
```

PL/TclU is the same, except that the language has to be specified as `pltclu`.

The body of the function is simply a piece of Tcl script. When the function is called, the argument values are passed to the Tcl script as variables named `1 ... n`. The result is returned from the Tcl code in the usual way, with a `return` statement. In a procedure, the return value from the Tcl code is ignored.

For example, a function returning the greater of two integer values could be defined as:

```
CREATE FUNCTION tcl_max(integer, integer) RETURNS integer AS $$
    if {$1 > $2} {return $1}
    return $2
$$ LANGUAGE pltcl STRICT;
```

Note the clause `STRICT`, which saves us from having to think about null input values: if a null value is passed, the function will not be called at all, but will just return a null result automatically.

In a nonstrict function, if the actual value of an argument is null, the corresponding `$n` variable will be set to an empty string. To detect whether a particular argument is null, use the function `argisnull`.

For example, suppose that we wanted `tcl_max` with one null and one nonnull argument to return the nonnull argument, rather than null:

```
CREATE FUNCTION tcl_max(integer, integer) RETURNS integer AS $$
  if {[argisnull 1]} {
    if {[argisnull 2]} { return_null }
    return $2
  }
  if {[argisnull 2]} { return $1 }
  if {$1 > $2} {return $1}
  return $2
$$ LANGUAGE pltcl;
```

As shown above, to return a null value from a PL/Tcl function, execute `return_null`. This can be done whether the function is strict or not.

Composite-type arguments are passed to the function as Tcl arrays. The element names of the array are the attribute names of the composite type. If an attribute in the passed row has the null value, it will not appear in the array. Here is an example:

```
CREATE TABLE employee (
  name text,
  salary integer,
  age integer
);

CREATE FUNCTION overpaid(employee) RETURNS boolean AS $$
  if {200000.0 < $1(salary)} {
    return "t"
  }
  if {$1(age) < 30 && 100000.0 < $1(salary)} {
    return "t"
  }
  return "f"
$$ LANGUAGE pltcl;
```

PL/Tcl functions can return composite-type results, too. To do this, the Tcl code must return a list of column name/value pairs matching the expected result type. Any column names omitted from the list are returned as nulls, and an error is raised if there are unexpected column names. Here is an example:

```
CREATE FUNCTION square_cube(in int, out squared int, out cubed int) AS $$
  return [list squared [expr {$1 * $1}] cubed [expr {$1 * $1 * $1}]]
$$ LANGUAGE pltcl;
```

Output arguments of procedures are returned in the same way, for example:

```
CREATE PROCEDURE tcl_triple(INOUT a integer, INOUT b integer) AS $$
  return [list a [expr {$1 * 3}] b [expr {$2 * 3}]]
$$ LANGUAGE pltcl;

CALL tcl_triple(5, 10);
```

Tip

The result list can be made from an array representation of the desired tuple with the `array get` Tcl command. For example:

```
CREATE FUNCTION raise_pay(employee, delta int) RETURNS employee AS $$
  set 1(salary) [expr {$1(salary) + $2}]
  return [array get 1]
$$ LANGUAGE pltcl;
```

PL/Tcl functions can return sets. To do this, the Tcl code should call `return_next` once per row to be returned, passing either the appropriate value when returning a scalar type, or a list of column name/value pairs when returning a composite type. Here is an example returning a scalar type:

```
CREATE FUNCTION sequence(int, int) RETURNS SETOF int AS $$
  for {set i $1} {$i < $2} {incr i} {
    return_next $i
  }
$$ LANGUAGE pltcl;
```

and here is one returning a composite type:

```
CREATE FUNCTION table_of_squares(int, int) RETURNS TABLE (x int, x2 int) AS $$
  for {set i $1} {$i < $2} {incr i} {
    return_next [list x $i x2 [expr {$i * $i}]]
  }
$$ LANGUAGE pltcl;
```

45.3. Data Values in PL/Tcl

The argument values supplied to a PL/Tcl function's code are simply the input arguments converted to text form (just as if they had been displayed by a `SELECT` statement). Conversely, the `return` and `return_next` commands will accept any string that is acceptable input format for the function's declared result type, or for the specified column of a composite result type.

45.4. Global Data in PL/Tcl

Sometimes it is useful to have some global data that is held between two calls to a function or is shared between different functions. This is easily done in PL/Tcl, but there are some restrictions that must be understood.

For security reasons, PL/Tcl executes functions called by any one SQL role in a separate Tcl interpreter for that role. This prevents accidental or malicious interference by one user with the behavior of another user's PL/Tcl functions. Each such interpreter will have its own values for any “global” Tcl variables. Thus, two PL/Tcl functions will share the same global variables if and only if they are executed by the same SQL role. In an application wherein a single session executes code under multiple SQL roles (via `SECURITY DEFINER` functions, use of `SET ROLE`, etc.) you may need to take explicit steps to ensure that PL/Tcl functions can share data. To do that, make sure that functions that should communicate are owned by the same user, and mark them `SECURITY DEFINER`. You must of course take care that such functions can't be used to do anything unintended.

All PL/TclU functions used in a session execute in the same Tcl interpreter, which of course is distinct from the interpreter(s) used for PL/Tcl functions. So global data is automatically shared between PL/TclU functions. This is not considered a security risk because all PL/TclU functions execute at the same trust level, namely that of a database superuser.

To help protect PL/Tcl functions from unintentionally interfering with each other, a global array is made available to each function via the `upvar` command. The global name of this variable is the function's internal name, and the local name is `GD`. It is recommended that `GD` be used for persistent private data of a function. Use regular Tcl global variables only for values that you specifically intend to be shared among multiple functions. (Note that the `GD` arrays are only global within a particular interpreter, so they do not bypass the security restrictions mentioned above.)

An example of using `GD` appears in the `spi_execp` example below.

45.5. Database Access from PL/Tcl

In this section, we follow the usual Tcl convention of using question marks, rather than brackets, to indicate an optional element in a syntax synopsis. The following commands are available to access the database from the body of a PL/Tcl function:

```
spi_exec ?-count n? ?-array name? command ?loop-body?
```

Executes an SQL command given as a string. An error in the command causes an error to be raised. Otherwise, the return value of `spi_exec` is the number of rows processed (selected, inserted, updated, or deleted) by the command, or zero if the command is a utility statement. In addition, if the command is a `SELECT` statement, the values of the selected columns are placed in Tcl variables as described below.

The optional `-count` value tells `spi_exec` to stop once *n* rows have been retrieved, much as if the query included a `LIMIT` clause. If *n* is zero, the query is run to completion, the same as when `-count` is omitted.

If the command is a `SELECT` statement, the values of the result columns are placed into Tcl variables named after the columns. If the `-array` option is given, the column values are instead stored into elements of the named associative array, with the column names used as array indexes. In addition, the current row number within the result (counting from zero) is stored into the array element named `".tupno"`, unless that name is in use as a column name in the result.

If the command is a `SELECT` statement and no *loop-body* script is given, then only the first row of results are stored into Tcl variables or array elements; remaining rows, if any, are ignored. No storing occurs if the query returns no rows. (This case can be detected by checking the result of `spi_exec`.) For example:

```
spi_exec "SELECT count(*) AS cnt FROM pg_proc"
```

will set the Tcl variable `$cnt` to the number of rows in the `pg_proc` system catalog.

If the optional *loop-body* argument is given, it is a piece of Tcl script that is executed once for each row in the query result. (*loop-body* is ignored if the given command is not a `SELECT`.) The values of the current row's columns are stored into Tcl variables or array elements before each iteration. For example:

```
spi_exec -array C "SELECT * FROM pg_class" {  
    elog DEBUG "have table $C(relname)"  
}
```

will print a log message for every row of `pg_class`. This feature works similarly to other Tcl looping constructs; in particular `continue` and `break` work in the usual way inside the loop body.

If a column of a query result is null, the target variable for it is “unset” rather than being set.

```
spi_prepare query typelist
```

Prepares and saves a query plan for later execution. The saved plan will be retained for the life of the current session.

The query can use parameters, that is, placeholders for values to be supplied whenever the plan is actually executed. In the query string, refer to parameters by the symbols `$1 ... $n`. If the query uses parameters, the names of the parameter types must be given as a Tcl list. (Write an empty list for *typelist* if no parameters are used.)

The return value from `spi_prepare` is a query ID to be used in subsequent calls to `spi_execp`. See `spi_execp` for an example.

```
spi_execp ?-count n? ?-array name? ?-nulls string? queryid ?value-list? ?loop-body?
```

Executes a query previously prepared with `spi_prepare`. *queryid* is the ID returned by `spi_prepare`. If the query references parameters, a *value-list* must be supplied. This is a Tcl list of actual values for the parameters. The list must be the same length as the parameter type list previously given to `spi_prepare`. Omit *value-list* if the query has no parameters.

The optional value for `-nulls` is a string of spaces and 'n' characters telling `spi_execp` which of the parameters are null values. If given, it must have exactly the same length as the *value-list*. If it is not given, all the parameter values are nonnull.

Except for the way in which the query and its parameters are specified, `spi_execp` works just like `spi_exec`. The `-count`, `-array`, and `loop-body` options are the same, and so is the result value.

Here's an example of a PL/Tcl function using a prepared plan:

```
CREATE FUNCTION t1_count(integer, integer) RETURNS integer AS $$
    if {![ info exists GD(plan) ]} {
        # prepare the saved plan on the first call
        set GD(plan) [ spi_prepare \
            "SELECT count(*) AS cnt FROM t1 WHERE num >= \ $1 AND num <= \ $2" \
            [ list int4 int4 ] ]
    }
    spi_execp -count 1 $GD(plan) [ list $1 $2 ]
    return $cnt
$$ LANGUAGE pltcl;
```

We need backslashes inside the query string given to `spi_prepare` to ensure that the `$n` markers will be passed through to `spi_prepare` as-is, and not replaced by Tcl variable substitution.

subtransaction command

The Tcl script contained in *command* is executed within an SQL subtransaction. If the script returns an error, that entire subtransaction is rolled back before returning the error out to the surrounding Tcl code. See [Section 45.9](#) for more details and an example.

quote string

Doubles all occurrences of single quote and backslash characters in the given string. This can be used to safely quote strings that are to be inserted into SQL commands given to `spi_exec` or `spi_prepare`. For example, think about an SQL command string like:

```
"SELECT '$val' AS ret"
```

where the Tcl variable `val` actually contains `doesn't`. This would result in the final command string:

```
SELECT 'doesn't' AS ret
```

which would cause a parse error during `spi_exec` or `spi_prepare`. To work properly, the submitted command should contain:

```
SELECT 'doesn''t' AS ret
```

which can be formed in PL/Tcl using:

```
"SELECT '[ quote $val ]' AS ret"
```

One advantage of `spi_execp` is that you don't have to quote parameter values like this, since the parameters are never parsed as part of an SQL command string.

elog level msg

Emits a log or error message. Possible levels are `DEBUG`, `LOG`, `INFO`, `NOTICE`, `WARNING`, `ERROR`, and `FATAL`. `ERROR` raises an error condition; if this is not trapped by the surrounding Tcl code, the error propagates out to the calling query, causing the current transaction or subtransaction to be aborted. This is effectively the same as the Tcl `error` command. `FATAL` aborts the transaction and causes the current session to shut down. (There is probably no good reason to use this error level in PL/Tcl functions, but it's provided for completeness.) The other levels only generate messages of different priority levels. Whether messages of a particular priority are reported to the client, written to the server log, or both is controlled by the [log_min_messages](#) and [client_min_messages](#) configuration variables. See [Chapter 19](#) and [Section 45.8](#) for more information.

45.6. Trigger Functions in PL/Tcl

Trigger functions can be written in PL/Tcl. Postgres Pro requires that a function that is to be called as a trigger must be declared as a function with no arguments and a return type of `trigger`.

The information from the trigger manager is passed to the function body in the following variables:

`$TG_name`

The name of the trigger from the `CREATE TRIGGER` statement.

`$TG_relid`

The object ID of the table that caused the trigger function to be invoked.

`$TG_table_name`

The name of the table that caused the trigger function to be invoked.

`$TG_table_schema`

The schema of the table that caused the trigger function to be invoked.

`$TG_relatts`

A Tcl list of the table column names, prefixed with an empty list element. So looking up a column name in the list with Tcl's `lsearch` command returns the element's number starting with 1 for the first column, the same way the columns are customarily numbered in Postgres Pro. (Empty list elements also appear in the positions of columns that have been dropped, so that the attribute numbering is correct for columns to their right.)

`$TG_when`

The string `BEFORE`, `AFTER`, or `INSTEAD OF`, depending on the type of trigger event.

`$TG_level`

The string `ROW` or `STATEMENT` depending on the type of trigger event.

`$TG_op`

The string `INSERT`, `UPDATE`, `DELETE`, or `TRUNCATE` depending on the type of trigger event.

`$NEW`

An associative array containing the values of the new table row for `INSERT` or `UPDATE` actions, or empty for `DELETE`. The array is indexed by column name. Columns that are null will not appear in the array. This is not set for statement-level triggers.

`$OLD`

An associative array containing the values of the old table row for `UPDATE` or `DELETE` actions, or empty for `INSERT`. The array is indexed by column name. Columns that are null will not appear in the array. This is not set for statement-level triggers.

`$args`

A Tcl list of the arguments to the function as given in the `CREATE TRIGGER` statement. These arguments are also accessible as `$1 ... $n` in the function body.

The return value from a trigger function can be one of the strings `OK` or `SKIP`, or a list of column name/value pairs. If the return value is `OK`, the operation (`INSERT/UPDATE/DELETE`) that fired the trigger will proceed normally. `SKIP` tells the trigger manager to silently suppress the operation for this row. If a list is returned, it tells PL/Tcl to return a modified row to the trigger manager; the contents of the modified row are specified by the column names and values in the list. Any columns not mentioned in the list are set to null. Returning a modified row is only meaningful for row-level `BEFORE INSERT` or `UPDATE` triggers, for which the modified row will be inserted instead of the one given in `$NEW`; or for row-level `INSTEAD OF INSERT` or `UPDATE` triggers where the returned row is used as the source data for `INSERT RETURNING` or `UPDATE RETURNING` clauses. In row-level `BEFORE DELETE` or `INSTEAD OF DELETE` triggers, returning

a modified row has the same effect as returning `OK`, that is the operation proceeds. The trigger return value is ignored for all other types of triggers.

Tip

The result list can be made from an array representation of the modified tuple with the `array get` Tcl command.

Here's a little example trigger function that forces an integer value in a table to keep track of the number of updates that are performed on the row. For new rows inserted, the value is initialized to 0 and then incremented on every update operation.

```
CREATE FUNCTION trigfunc_modcount() RETURNS trigger AS $$
    switch $TG_op {
        INSERT {
            set NEW($1) 0
        }
        UPDATE {
            set NEW($1) $OLD($1)
            incr NEW($1)
        }
        default {
            return OK
        }
    }
    return [array get NEW]
$$ LANGUAGE pltcl;

CREATE TABLE mytab (num integer, description text, modcnt integer);

CREATE TRIGGER trig_mytab_modcount BEFORE INSERT OR UPDATE ON mytab
    FOR EACH ROW EXECUTE FUNCTION trigfunc_modcount('modcnt');
```

Notice that the trigger function itself does not know the column name; that's supplied from the trigger arguments. This lets the trigger function be reused with different tables.

45.7. Event Trigger Functions in PL/Tcl

Event trigger functions can be written in PL/Tcl. Postgres Pro requires that a function that is to be called as an event trigger must be declared as a function with no arguments and a return type of `event_trigger`.

The information from the trigger manager is passed to the function body in the following variables:

`$TG_event`

The name of the event the trigger is fired for.

`$TG_tag`

The command tag for which the trigger is fired.

The return value of the trigger function is ignored.

Here's a little example event trigger function that simply raises a `NOTICE` message each time a supported command is executed:

```
CREATE OR REPLACE FUNCTION tclsnitch() RETURNS event_trigger AS $$
    elog NOTICE "tclsnitch: $TG_event $TG_tag"
$$ LANGUAGE pltcl;
```

```
CREATE EVENT TRIGGER tcl_a_snitch ON ddl_command_start EXECUTE FUNCTION tclsnitch();
```

45.8. Error Handling in PL/Tcl

Tcl code within or called from a PL/Tcl function can raise an error, either by executing some invalid operation or by generating an error using the Tcl `error` command or PL/Tcl's `elog` command. Such errors can be caught within Tcl using the Tcl `catch` command. If an error is not caught but is allowed to propagate out to the top level of execution of the PL/Tcl function, it is reported as an SQL error in the function's calling query.

Conversely, SQL errors that occur within PL/Tcl's `spi_exec`, `spi_prepare`, and `spi_execp` commands are reported as Tcl errors, so they are catchable by Tcl's `catch` command. (Each of these PL/Tcl commands runs its SQL operation in a subtransaction, which is rolled back on error, so that any partially-completed operation is automatically cleaned up.) Again, if an error propagates out to the top level without being caught, it turns back into an SQL error.

Tcl provides an `errorCode` variable that can represent additional information about an error in a form that is easy for Tcl programs to interpret. The contents are in Tcl list format, and the first word identifies the subsystem or library reporting the error; beyond that the contents are left to the individual subsystem or library. For database errors reported by PL/Tcl commands, the first word is `POSTGRES`, the second word is the Postgres Pro version number, and additional words are field name/value pairs providing detailed information about the error. Fields `SQLSTATE`, `condition`, and `message` are always supplied (the first two represent the error code and condition name as shown in [Appendix A](#)). Fields that may be present include `detail`, `hint`, `context`, `schema`, `table`, `column`, `datatype`, `constraint`, `statement`, `cursor_position`, `filename`, `lineno`, and `funcname`.

A convenient way to work with PL/Tcl's `errorCode` information is to load it into an array, so that the field names become array subscripts. Code for doing that might look like

```
if {[catch { spi_exec $sql_command }]} {
    if {[lindex $::errorCode 0] == "POSTGRES"} {
        array set errorArray $::errorCode
        if {$errorArray(condition) == "undefined_table"} {
            # deal with missing table
        } else {
            # deal with some other type of SQL error
        }
    }
}
```

(The double colons explicitly specify that `errorCode` is a global variable.)

45.9. Explicit Subtransactions in PL/Tcl

Recovering from errors caused by database access as described in [Section 45.8](#) can lead to an undesirable situation where some operations succeed before one of them fails, and after recovering from that error the data is left in an inconsistent state. PL/Tcl offers a solution to this problem in the form of explicit subtransactions.

Consider a function that implements a transfer between two accounts:

```
CREATE FUNCTION transfer_funds() RETURNS void AS $$
    if [catch {
        spi_exec "UPDATE accounts SET balance = balance - 100 WHERE account_name =
'joe'"
        spi_exec "UPDATE accounts SET balance = balance + 100 WHERE account_name =
'mary'"
    } errmsg] {
        set result [format "error transferring funds: %s" $errmsg]
```

```
    } else {
        set result "funds transferred successfully"
    }
    spi_exec "INSERT INTO operations (result) VALUES ('[quote $result]')"
$$ LANGUAGE pltcl;
```

If the second `UPDATE` statement results in an exception being raised, this function will log the failure, but the result of the first `UPDATE` will nevertheless be committed. In other words, the funds will be withdrawn from Joe's account, but will not be transferred to Mary's account. This happens because each `spi_exec` is a separate subtransaction, and only one of those subtransactions got rolled back.

To handle such cases, you can wrap multiple database operations in an explicit subtransaction, which will succeed or roll back as a whole. PL/Tcl provides a `subtransaction` command to manage this. We can rewrite our function as:

```
CREATE FUNCTION transfer_funds2() RETURNS void AS $$
    if [catch {
        subtransaction {
            spi_exec "UPDATE accounts SET balance = balance - 100 WHERE account_name =
'joe'"
            spi_exec "UPDATE accounts SET balance = balance + 100 WHERE account_name =
'mary'"
        }
    } errmsg] {
        set result [format "error transferring funds: %s" $errmsg]
    } else {
        set result "funds transferred successfully"
    }
    spi_exec "INSERT INTO operations (result) VALUES ('[quote $result]')"
$$ LANGUAGE pltcl;
```

Note that use of `catch` is still required for this purpose. Otherwise the error would propagate to the top level of the function, preventing the desired insertion into the `operations` table. The `subtransaction` command does not trap errors, it only assures that all database operations executed inside its scope will be rolled back together when an error is reported.

A rollback of an explicit subtransaction occurs on any error reported by the contained Tcl code, not only errors originating from database access. Thus a regular Tcl exception raised inside a `subtransaction` command will also cause the subtransaction to be rolled back. However, non-error exits out of the contained Tcl code (for instance, due to `return`) do not cause a rollback.

45.10. Transaction Management

In a procedure called from the top level or an anonymous code block (`DO` command) called from the top level it is possible to control transactions. To commit the current transaction, call the `commit` command. To roll back the current transaction, call the `rollback` command. (Note that it is not possible to run the SQL commands `COMMIT` or `ROLLBACK` via `spi_exec` or similar. It has to be done using these functions.) After a transaction is ended, a new transaction is automatically started, so there is no separate command for that.

Here is an example:

```
CREATE PROCEDURE transaction_test1()
LANGUAGE pltcl
AS $$
for {set i 0} {$i < 10} {incr i} {
    spi_exec "INSERT INTO test1 (a) VALUES ($i)"
    if {$i % 2 == 0} {
        commit
    } else {
```

```
        rollback
    }
}
$$;
```

```
CALL transaction_test1();
```

Transactions cannot be ended when an explicit subtransaction is active.

45.11. PL/Tcl Configuration

This section lists configuration parameters that affect PL/Tcl.

```
pltcl.start_proc (string)
```

This parameter, if set to a nonempty string, specifies the name (possibly schema-qualified) of a parameterless PL/Tcl function that is to be executed whenever a new Tcl interpreter is created for PL/Tcl. Such a function can perform per-session initialization, such as loading additional Tcl code. A new Tcl interpreter is created when a PL/Tcl function is first executed in a database session, or when an additional interpreter has to be created because a PL/Tcl function is called by a new SQL role.

The referenced function must be written in the `pltcl` language, and must not be marked `SECURITY DEFINER`. (These restrictions ensure that it runs in the interpreter it's supposed to initialize.) The current user must have permission to call it, too.

If the function fails with an error it will abort the function call that caused the new interpreter to be created and propagate out to the calling query, causing the current transaction or subtransaction to be aborted. Any actions already done within Tcl won't be undone; however, that interpreter won't be used again. If the language is used again the initialization will be attempted again within a fresh Tcl interpreter.

Only superusers can change this setting. Although this setting can be changed within a session, such changes will not affect Tcl interpreters that have already been created.

```
pltclu.start_proc (string)
```

This parameter is exactly like `pltcl.start_proc`, except that it applies to PL/TclU. The referenced function must be written in the `pltclu` language.

45.12. Tcl Procedure Names

In Postgres Pro, the same function name can be used for different function definitions as long as the number of arguments or their types differ. Tcl, however, requires all procedure names to be distinct. PL/Tcl deals with this by making the internal Tcl procedure names contain the object ID of the function from the system table `pg_proc` as part of their name. Thus, Postgres Pro functions with the same name and different argument types will be different Tcl procedures, too. This is not normally a concern for a PL/Tcl programmer, but it might be visible when debugging.

Chapter 46. PL/Perl — Perl Procedural Language

PL/Perl is a loadable procedural language that enables you to write Postgres Pro functions and procedures in the [Perl programming language](#).

The main advantage to using PL/Perl is that this allows use, within stored functions and procedures, of the manyfold “string munging” operators and functions available for Perl. Parsing complex strings might be easier using Perl than it is with the string functions and control structures provided in PL/pgSQL.

To install PL/Perl in a particular database, use `CREATE EXTENSION plperl`.

Tip

If a language is installed into `template1`, all subsequently created databases will have the language installed automatically.

Note

Users of source packages must specially enable the build of PL/Perl during the installation process. Users of binary packages might find PL/Perl in a separate subpackage.

46.1. PL/Perl Functions and Arguments

To create a function in the PL/Perl language, use the standard [CREATE FUNCTION](#) syntax:

```
CREATE FUNCTION funcname (argument-types)
RETURNS return-type
-- function attributes can go here
AS $$
    # PL/Perl function body goes here
$$ LANGUAGE plperl;
```

The body of the function is ordinary Perl code. In fact, the PL/Perl glue code wraps it inside a Perl subroutine. A PL/Perl function is called in a scalar context, so it can't return a list. You can return non-scalar values (arrays, records, and sets) by returning a reference, as discussed below.

In a PL/Perl procedure, any return value from the Perl code is ignored.

PL/Perl also supports anonymous code blocks called with the [DO](#) statement:

```
DO $$
    # PL/Perl code
$$ LANGUAGE plperl;
```

An anonymous code block receives no arguments, and whatever value it might return is discarded. Otherwise it behaves just like a function.

Note

The use of named nested subroutines is dangerous in Perl, especially if they refer to lexical variables in the enclosing scope. Because a PL/Perl function is wrapped in a subroutine, any named subroutine you place inside one will be nested. In general, it is far safer to create anonymous subroutines which you call via a coderef. For more information, see the entries for `Variable "%s"`

will not stay shared and Variable "%s" is not available in the `perldiag` man page, or search the Internet for “perl nested named subroutine”.

The syntax of the `CREATE FUNCTION` command requires the function body to be written as a string constant. It is usually most convenient to use dollar quoting (see [Section 4.1.2.4](#)) for the string constant. If you choose to use escape string syntax `E' '`, you must double any single quote marks (`'`) and backslashes (`\`) used in the body of the function (see [Section 4.1.2.1](#)).

Arguments and results are handled as in any other Perl subroutine: arguments are passed in `@_`, and a result value is returned with `return` or as the last expression evaluated in the function.

For example, a function returning the greater of two integer values could be defined as:

```
CREATE FUNCTION perl_max (integer, integer) RETURNS integer AS $$
    if ($_[0] > $_[1]) { return $_[0]; }
    return $_[1];
$$ LANGUAGE plperl;
```

Note

Arguments will be converted from the database's encoding to UTF-8 for use inside PL/Perl, and then converted from UTF-8 back to the database encoding upon return.

If an SQL null value is passed to a function, the argument value will appear as “undefined” in Perl. The above function definition will not behave very nicely with null inputs (in fact, it will act as though they are zeroes). We could add `STRICT` to the function definition to make Postgres Pro do something more reasonable: if a null value is passed, the function will not be called at all, but will just return a null result automatically. Alternatively, we could check for undefined inputs in the function body. For example, suppose that we wanted `perl_max` with one null and one nonnull argument to return the nonnull argument, rather than a null value:

```
CREATE FUNCTION perl_max (integer, integer) RETURNS integer AS $$
    my ($x, $y) = @_;
    if (not defined $x) {
        return undef if not defined $y;
        return $y;
    }
    return $x if not defined $y;
    return $x if $x > $y;
    return $y;
$$ LANGUAGE plperl;
```

As shown above, to return an SQL null value from a PL/Perl function, return an undefined value. This can be done whether the function is strict or not.

Anything in a function argument that is not a reference is a string, which is in the standard Postgres Pro external text representation for the relevant data type. In the case of ordinary numeric or text types, Perl will just do the right thing and the programmer will normally not have to worry about it. However, in other cases the argument will need to be converted into a form that is more usable in Perl. For example, the `decode_bytea` function can be used to convert an argument of type `bytea` into unescaped binary.

Similarly, values passed back to Postgres Pro must be in the external text representation format. For example, the `encode_bytea` function can be used to escape binary data for a return value of type `bytea`.

One case that is particularly important is boolean values. As just stated, the default behavior for `bool` values is that they are passed to Perl as text, thus either `'t'` or `'f'`. This is problematic, since Perl will not treat `'f'` as false! It is possible to improve matters by using a “transform” (see [CREATE TRANSFORM](#)). Suitable transforms are provided by the `bool_plperl` extension. To use it, install the extension:

```
CREATE EXTENSION bool_plperl; -- or bool_plperl_u for PL/PerlU
```

Then use the `TRANSFORM` function attribute for a PL/Perl function that takes or returns `bool`, for example:

```
CREATE FUNCTION perl_and(bool, bool) RETURNS bool
TRANSFORM FOR TYPE bool
AS $$
    my ($a, $b) = @_;
    return $a && $b;
$$ LANGUAGE plperl;
```

When this transform is applied, `bool` arguments will be seen by Perl as being 1 or empty, thus properly true or false. If the function result is type `bool`, it will be true or false according to whether Perl would evaluate the returned value as true. Similar transformations are also performed for boolean query arguments and results of SPI queries performed inside the function ([Section 46.3.1](#)).

Perl can return PostgreSQL arrays as references to Perl arrays. Here is an example:

```
CREATE OR REPLACE FUNCTION returns_array()
RETURNS text[][] AS $$
    return [['a"b', 'c,d'], ['e\\f', 'g']];
$$ LANGUAGE plperl;

select returns_array();
```

Perl passes PostgreSQL arrays as a blessed `PostgreSQL::InServer::ARRAY` object. This object may be treated as an array reference or a string, allowing for backward compatibility with Perl code written for PostgreSQL versions below 9.1 to run. For example:

```
CREATE OR REPLACE FUNCTION concat_array_elements(text[]) RETURNS TEXT AS $$
    my $arg = shift;
    my $result = "";
    return undef if (!defined $arg);

    # as an array reference
    for (@$arg) {
        $result .= $_;
    }

    # also works as a string
    $result .= $arg;

    return $result;
$$ LANGUAGE plperl;

SELECT concat_array_elements(ARRAY['PL', '/', 'Perl']);
```

Note

Multidimensional arrays are represented as references to lower-dimensional arrays of references in a way common to every Perl programmer.

Composite-type arguments are passed to the function as references to hashes. The keys of the hash are the attribute names of the composite type. Here is an example:

```
CREATE TABLE employee (
    name text,
    basesalary integer,
    bonus integer
);
```



```
CREATE FUNCTION empcomp(employee) RETURNS integer AS $$
    my ($emp) = @_;
    return $emp->{basesalary} + $emp->{bonus};
$$ LANGUAGE plperl;
```

```
SELECT name, empcomp(employee.*) FROM employee;
```

A PL/Perl function can return a composite-type result using the same approach: return a reference to a hash that has the required attributes. For example:

```
CREATE TYPE testrowperl AS (f1 integer, f2 text, f3 text);

CREATE OR REPLACE FUNCTION perl_row() RETURNS testrowperl AS $$
    return {f2 => 'hello', f1 => 1, f3 => 'world'};
$$ LANGUAGE plperl;

SELECT * FROM perl_row();
```

Any columns in the declared result data type that are not present in the hash will be returned as null values.

Similarly, output arguments of procedures can be returned as a hash reference:

```
CREATE PROCEDURE perl_triple(INOUT a integer, INOUT b integer) AS $$
    my ($a, $b) = @_;
    return {a => $a * 3, b => $b * 3};
$$ LANGUAGE plperl;

CALL perl_triple(5, 10);
```

PL/Perl functions can also return sets of either scalar or composite types. Usually you'll want to return rows one at a time, both to speed up startup time and to keep from queuing up the entire result set in memory. You can do this with `return_next` as illustrated below. Note that after the last `return_next`, you must put either `return` or (better) `return undef`.

```
CREATE OR REPLACE FUNCTION perl_set_int(int)
RETURNS SETOF INTEGER AS $$
    foreach (0..$_[0]) {
        return_next($_);
    }
    return undef;
$$ LANGUAGE plperl;

SELECT * FROM perl_set_int(5);

CREATE OR REPLACE FUNCTION perl_set()
RETURNS SETOF testrowperl AS $$
    return_next({ f1 => 1, f2 => 'Hello', f3 => 'World' });
    return_next({ f1 => 2, f2 => 'Hello', f3 => 'PostgreSQL' });
    return_next({ f1 => 3, f2 => 'Hello', f3 => 'PL/Perl' });
    return undef;
$$ LANGUAGE plperl;
```

For small result sets, you can return a reference to an array that contains either scalars, references to arrays, or references to hashes for simple types, array types, and composite types, respectively. Here are some simple examples of returning the entire result set as an array reference:

```
CREATE OR REPLACE FUNCTION perl_set_int(int) RETURNS SETOF INTEGER AS $$
    return [0..$_[0]];
$$ LANGUAGE plperl;
```

```
SELECT * FROM perl_set_int(5);

CREATE OR REPLACE FUNCTION perl_set() RETURNS SETOF testrowperl AS $$
    return [
        { f1 => 1, f2 => 'Hello', f3 => 'World' },
        { f1 => 2, f2 => 'Hello', f3 => 'PostgreSQL' },
        { f1 => 3, f2 => 'Hello', f3 => 'PL/Perl' }
    ];
$$ LANGUAGE plperl;

SELECT * FROM perl_set();
```

If you wish to use the `strict` pragma with your code you have a few options. For temporary global use you can `SET plperl.use_strict` to true. This will affect subsequent compilations of PL/Perl functions, but not functions already compiled in the current session. For permanent global use you can set `plperl.use_strict` to true in the `postgresql.conf` file.

For permanent use in specific functions you can simply put:

```
use strict;

at the top of the function body.
```

The `feature` pragma is also available to use if your Perl is version 5.10.0 or higher.

46.2. Data Values in PL/Perl

The argument values supplied to a PL/Perl function's code are simply the input arguments converted to text form (just as if they had been displayed by a `SELECT` statement). Conversely, the `return` and `return_next` commands will accept any string that is acceptable input format for the function's declared return type.

If this behavior is inconvenient for a particular case, it can be improved by using a transform, as already illustrated for `bool` values. Several examples of transform modules are included in the PostgreSQL distribution.

46.3. Built-in Functions

46.3.1. Database Access from PL/Perl

Access to the database itself from your Perl function can be done via the following functions:

```
spi_exec_query(query [, limit])
```

`spi_exec_query` executes an SQL command and returns the entire row set as a reference to an array of hash references. If `limit` is specified and is greater than zero, then `spi_exec_query` retrieves at most `limit` rows, much as if the query included a `LIMIT` clause. Omitting `limit` or specifying it as zero results in no row limit.

You should only use this command when you know that the result set will be relatively small. Here is an example of a query (`SELECT` command) with the optional maximum number of rows:

```
$rv = spi_exec_query('SELECT * FROM my_table', 5);
```

This returns up to 5 rows from the table `my_table`. If `my_table` has a column `my_column`, you can get that value from row `$i` of the result like this:

```
$foo = $rv->{rows}[$i]->{my_column};
```

The total number of rows returned from a `SELECT` query can be accessed like this:

```
$nrows = $rv->{processed}
```

Here is an example using a different command type:

```
$query = "INSERT INTO my_table VALUES (1, 'test')";
$rv = spi_exec_query($query);
```

You can then access the command status (e.g., SPI_OK_INSERT) like this:

```
$res = $rv->{status};
```

To get the number of rows affected, do:

```
$nrows = $rv->{processed};
```

Here is a complete example:

```
CREATE TABLE test (
    i int,
    v varchar
);

INSERT INTO test (i, v) VALUES (1, 'first line');
INSERT INTO test (i, v) VALUES (2, 'second line');
INSERT INTO test (i, v) VALUES (3, 'third line');
INSERT INTO test (i, v) VALUES (4, 'immortal');

CREATE OR REPLACE FUNCTION test_munge() RETURNS SETOF test AS $$
    my $rv = spi_exec_query('select i, v from test;');
    my $status = $rv->{status};
    my $nrows = $rv->{processed};
    foreach my $rn (0 .. $nrows - 1) {
        my $row = $rv->{rows}[$rn];
        $row->{i} += 200 if defined($row->{i});
        $row->{v} =~ tr/A-Za-z/a-zA-Z/ if (defined($row->{v}));
        return_next($row);
    }
    return undef;
$$ LANGUAGE plperl;

SELECT * FROM test_munge();

spi_query(command)
spi_fetchrow(cursor)
spi_cursor_close(cursor)
```

spi_query and **spi_fetchrow** work together as a pair for row sets which might be large, or for cases where you wish to return rows as they arrive. **spi_fetchrow** works *only* with **spi_query**. The following example illustrates how you use them together:

```
CREATE TYPE foo_type AS (the_num INTEGER, the_text TEXT);

CREATE OR REPLACE FUNCTION lotsa_md5 (INTEGER) RETURNS SETOF foo_type AS $$
    use Digest::MD5 qw(md5_hex);
    my $file = '/usr/share/dict/words';
    my $t = localtime;
    elog(NOTICE, "opening file $file at $t");
    open my $fh, '<', $file # ooh, it's a file access!
        or elog(ERROR, "cannot open $file for reading: $!");
    my @words = <$fh>;
    close $fh;
    $t = localtime;
    elog(NOTICE, "closed file $file at $t");
    chomp(@words);
    my $row;
    my $sth = spi_query("SELECT * FROM generate_series(1,$_[0]) AS b(a)");
```

```
while (defined ($row = spi_fetchrow($sth))) {
    return_next({
        the_num => $row->{a},
        the_text => md5_hex($words[rand @words])
    });
}
return;
$$ LANGUAGE plperl;

SELECT * from lotsa_md5(500);
```

Normally, `spi_fetchrow` should be repeated until it returns `undef`, indicating that there are no more rows to read. The cursor returned by `spi_query` is automatically freed when `spi_fetchrow` returns `undef`. If you do not wish to read all the rows, instead call `spi_cursor_close` to free the cursor. Failure to do so will result in memory leaks.

```
spi_prepare(command, argument types)
spi_query_prepared(plan, arguments)
spi_exec_prepared(plan [, attributes], arguments)
spi_freeplan(plan)
```

`spi_prepare`, `spi_query_prepared`, `spi_exec_prepared`, and `spi_freeplan` implement the same functionality but for prepared queries. `spi_prepare` accepts a query string with numbered argument placeholders (`$1`, `$2`, etc.) and a string list of argument types:

```
$plan = spi_prepare('SELECT * FROM test WHERE id > $1 AND name = $2',
                    'INTEGER', 'TEXT');
```

Once a query plan is prepared by a call to `spi_prepare`, the plan can be used instead of the string query, either in `spi_exec_prepared`, where the result is the same as returned by `spi_exec_query`, or in `spi_query_prepared` which returns a cursor exactly as `spi_query` does, which can be later passed to `spi_fetchrow`. The optional second parameter to `spi_exec_prepared` is a hash reference of attributes; the only attribute currently supported is `limit`, which sets the maximum number of rows returned from the query. Omitting `limit` or specifying it as zero results in no row limit.

The advantage of prepared queries is that it is possible to use one prepared plan for more than one query execution. After the plan is not needed anymore, it can be freed with `spi_freeplan`:

```
CREATE OR REPLACE FUNCTION init() RETURNS VOID AS $$
    $_SHARED{my_plan} = spi_prepare('SELECT (now() + $1)::date AS now',
                                    'INTERVAL');
$$ LANGUAGE plperl;

CREATE OR REPLACE FUNCTION add_time( INTERVAL ) RETURNS TEXT AS $$
    return spi_exec_prepared(
        $_SHARED{my_plan},
        $_[0]
    )->{rows}->[0]->{now};
$$ LANGUAGE plperl;

CREATE OR REPLACE FUNCTION done() RETURNS VOID AS $$
    spi_freeplan( $_SHARED{my_plan});
    undef $_SHARED{my_plan};
$$ LANGUAGE plperl;

SELECT init();
SELECT add_time('1 day'), add_time('2 days'), add_time('3 days');
SELECT done();
```

```
add_time | add_time | add_time
```

```
-----+-----+-----
2005-12-10 | 2005-12-11 | 2005-12-12
```

Note that the parameter subscript in `spi_prepare` is defined via `$1`, `$2`, `$3`, etc., so avoid declaring query strings in double quotes that might easily lead to hard-to-catch bugs.

Another example illustrates usage of an optional parameter in `spi_exec_prepared`:

```
CREATE TABLE hosts AS SELECT id, ('192.168.1.'||id)::inet AS address
                        FROM generate_series(1,3) AS id;

CREATE OR REPLACE FUNCTION init_hosts_query() RETURNS VOID AS $$
    $_SHARED{plan} = spi_prepare('SELECT * FROM hosts
                                WHERE address << $1', 'inet');
$$ LANGUAGE plperl;

CREATE OR REPLACE FUNCTION query_hosts(inet) RETURNS SETOF hosts AS $$
    return spi_exec_prepared(
        $_SHARED{plan},
        {limit => 2},
        $_[0]
    )->{rows};
$$ LANGUAGE plperl;

CREATE OR REPLACE FUNCTION release_hosts_query() RETURNS VOID AS $$
    spi_freeplan($_SHARED{plan});
    undef $_SHARED{plan};
$$ LANGUAGE plperl;

SELECT init_hosts_query();
SELECT query_hosts('192.168.1.0/30');
SELECT release_hosts_query();
```

```
      query_hosts
-----
(1,192.168.1.1)
(2,192.168.1.2)
(2 rows)
```

```
spi_commit()
spi_rollback()
```

Commit or roll back the current transaction. This can only be called in a procedure or anonymous code block (DO command) called from the top level. (Note that it is not possible to run the SQL commands COMMIT or ROLLBACK via `spi_exec_query` or similar. It has to be done using these functions.) After a transaction is ended, a new transaction is automatically started, so there is no separate function for that.

Here is an example:

```
CREATE PROCEDURE transaction_test1()
LANGUAGE plperl
AS $$
foreach my $i (0..9) {
    spi_exec_query("INSERT INTO test1 (a) VALUES ($i)");
    if ($i % 2 == 0) {
        spi_commit();
    } else {
        spi_rollback();
    }
}
}
```

```
$$;
```

```
CALL transaction_test1();
```

46.3.2. Utility Functions in PL/Perl

`elog(level, msg)`

Emit a log or error message. Possible levels are `DEBUG`, `LOG`, `INFO`, `NOTICE`, `WARNING`, and `ERROR`. `ERROR` raises an error condition; if this is not trapped by the surrounding Perl code, the error propagates out to the calling query, causing the current transaction or subtransaction to be aborted. This is effectively the same as the Perl `die` command. The other levels only generate messages of different priority levels. Whether messages of a particular priority are reported to the client, written to the server log, or both is controlled by the [log_min_messages](#) and [client_min_messages](#) configuration variables. See [Chapter 19](#) for more information.

`quote_literal(string)`

Return the given string suitably quoted to be used as a string literal in an SQL statement string. Embedded single-quotes and backslashes are properly doubled. Note that `quote_literal` returns `undef` on `undef` input; if the argument might be `undef`, `quote_nullable` is often more suitable.

`quote_nullable(string)`

Return the given string suitably quoted to be used as a string literal in an SQL statement string; or, if the argument is `undef`, return the unquoted string "NULL". Embedded single-quotes and backslashes are properly doubled.

`quote_ident(string)`

Return the given string suitably quoted to be used as an identifier in an SQL statement string. Quotes are added only if necessary (i.e., if the string contains non-identifier characters or would be case-folded). Embedded quotes are properly doubled.

`decode_bytea(string)`

Return the unescaped binary data represented by the contents of the given string, which should be `bytea` encoded.

`encode_bytea(string)`

Return the `bytea` encoded form of the binary data contents of the given string.

`encode_array_literal(array)`

`encode_array_literal(array, delimiter)`

Returns the contents of the referenced array as a string in array literal format (see [Section 8.15.2](#)). Returns the argument value unaltered if it's not a reference to an array. The delimiter used between elements of the array literal defaults to `,` if a delimiter is not specified or is `undef`.

`encode_typed_literal(value, typename)`

Converts a Perl variable to the value of the data type passed as a second argument and returns a string representation of this value. Correctly handles nested arrays and values of composite types.

`encode_array_constructor(array)`

Returns the contents of the referenced array as a string in array constructor format (see [Section 4.2.12](#)). Individual values are quoted using `quote_nullable`. Returns the argument value, quoted using `quote_nullable`, if it's not a reference to an array.

`looks_like_number(string)`

Returns a true value if the content of the given string looks like a number, according to Perl, returns false otherwise. Returns `undef` if the argument is `undef`. Leading and trailing space is ignored. `Inf` and `Infinity` are regarded as numbers.

```
is_array_ref(argument)
```

Returns a true value if the given argument may be treated as an array reference, that is, if ref of the argument is ARRAY or PostgreSQL::InServer::ARRAY. Returns false otherwise.

46.4. Global Values in PL/Perl

You can use the global hash %_SHARED to store data, including code references, between function calls for the lifetime of the current session.

Here is a simple example for shared data:

```
CREATE OR REPLACE FUNCTION set_var(name text, val text) RETURNS text AS $$
    if ($_SHARED{$_[0]} = $_[1]) {
        return 'ok';
    } else {
        return "cannot set shared variable $_[0] to $_[1]";
    }
$$ LANGUAGE plperl;
```

```
CREATE OR REPLACE FUNCTION get_var(name text) RETURNS text AS $$
    return $_SHARED{$_[0]};
$$ LANGUAGE plperl;
```

```
SELECT set_var('sample', 'Hello, PL/Perl!  How's tricks?');
SELECT get_var('sample');
```

Here is a slightly more complicated example using a code reference:

```
CREATE OR REPLACE FUNCTION myfuncs() RETURNS void AS $$
    $_SHARED{myquote} = sub {
        my $arg = shift;
        $arg =~ s/(['\\])/\\$1/g;
        return "'$arg'";
    };
$$ LANGUAGE plperl;

SELECT myfuncs(); /* initializes the function */

/* Set up a function that uses the quote function */

CREATE OR REPLACE FUNCTION use_quote(TEXT) RETURNS text AS $$
    my $text_to_quote = shift;
    my $qfunc = $_SHARED{myquote};
    return &$qfunc($text_to_quote);
$$ LANGUAGE plperl;
```

(You could have replaced the above with the one-liner `return $_SHARED{myquote}->($_[0]);` at the expense of readability.)

For security reasons, PL/Perl executes functions called by any one SQL role in a separate Perl interpreter for that role. This prevents accidental or malicious interference by one user with the behavior of another user's PL/Perl functions. Each such interpreter has its own value of the %_SHARED variable and other global state. Thus, two PL/Perl functions will share the same value of %_SHARED if and only if they are executed by the same SQL role. In an application wherein a single session executes code under multiple SQL roles (via SECURITY DEFINER functions, use of SET ROLE, etc.) you may need to take explicit steps to ensure that PL/Perl functions can share data via %_SHARED. To do that, make sure that functions that should communicate are owned by the same user, and mark them SECURITY DEFINER. You must of course take care that such functions can't be used to do anything unintended.

46.5. Trusted and Untrusted PL/Perl

Normally, PL/Perl is installed as a “trusted” programming language named `plperl`. In this setup, certain Perl operations are disabled to preserve security. In general, the operations that are restricted are those that interact with the environment. This includes file handle operations, `require`, and `use` (for external modules). There is no way to access internals of the database server process or to gain OS-level access with the permissions of the server process, as a C function can do. Thus, any unprivileged database user can be permitted to use this language.

Warning

Trusted PL/Perl relies on the Perl `Opcode` module to preserve security. Perl [documents](#) that the module is not effective for the trusted PL/Perl use case. If your security needs are incompatible with the uncertainty in that warning, consider executing `REVOKE USAGE ON LANGUAGE plperl FROM PUBLIC`.

Here is an example of a function that will not work because file system operations are not allowed for security reasons:

```
CREATE FUNCTION badfunc() RETURNS integer AS $$
    my $tmpfile = "/tmp/badfile";
    open my $fh, '>', $tmpfile
        or elog(ERROR, qq{could not open the file "$tmpfile": $!});
    print $fh "Testing writing to a file\n";
    close $fh or elog(ERROR, qq{could not close the file "$tmpfile": $!});
    return 1;
$$ LANGUAGE plperl;
```

The creation of this function will fail as its use of a forbidden operation will be caught by the validator.

Sometimes it is desirable to write Perl functions that are not restricted. For example, one might want a Perl function that sends mail. To handle these cases, PL/Perl can also be installed as an “untrusted” language (usually called PL/PerlU). In this case the full Perl language is available. When installing the language, the language name `plperlU` will select the untrusted PL/Perl variant.

The writer of a PL/PerlU function must take care that the function cannot be used to do anything unwanted, since it will be able to do anything that could be done by a user logged in as the database administrator. Note that the database system allows only database superusers to create functions in untrusted languages.

If the above function was created by a superuser using the language `plperlU`, execution would succeed.

In the same way, anonymous code blocks written in Perl can use restricted operations if the language is specified as `plperlU` rather than `plperl`, but the caller must be a superuser.

Note

While PL/Perl functions run in a separate Perl interpreter for each SQL role, all PL/PerlU functions executed in a given session run in a single Perl interpreter (which is not any of the ones used for PL/Perl functions). This allows PL/PerlU functions to share data freely, but no communication can occur between PL/Perl and PL/PerlU functions.

Note

Perl cannot support multiple interpreters within one process unless it was built with the appropriate flags, namely either `usemultiplicity` or `useithreads`. (`usemultiplicity` is preferred unless you actually need to use threads. For more details, see the `perlembed` man page.) If PL/Perl is used with a copy of Perl that was not built this way, then it is only possible to have one Perl

interpreter per session, and so any one session can only execute either PL/PerlU functions, or PL/Perl functions that are all called by the same SQL role.

46.6. PL/Perl Triggers

PL/Perl can be used to write trigger functions. In a trigger function, the hash reference `$_TD` contains information about the current trigger event. `$_TD` is a global variable, which gets a separate local value for each invocation of the trigger. The fields of the `$_TD` hash reference are:

`$_TD->{new}{foo}`

NEW value of column `foo`

`$_TD->{old}{foo}`

OLD value of column `foo`

`$_TD->{name}`

Name of the trigger being called

`$_TD->{event}`

Trigger event: INSERT, UPDATE, DELETE, TRUNCATE, or UNKNOWN

`$_TD->{when}`

When the trigger was called: BEFORE, AFTER, INSTEAD OF, or UNKNOWN

`$_TD->{level}`

The trigger level: ROW, STATEMENT, or UNKNOWN

`$_TD->{relid}`

OID of the table on which the trigger fired

`$_TD->{table_name}`

Name of the table on which the trigger fired

`$_TD->{relname}`

Name of the table on which the trigger fired. This has been deprecated, and could be removed in a future release. Please use `$_TD->{table_name}` instead.

`$_TD->{table_schema}`

Name of the schema in which the table on which the trigger fired, is

`$_TD->{argc}`

Number of arguments of the trigger function

`@{$_TD->{args}}`

Arguments of the trigger function. Does not exist if `$_TD->{argc}` is 0.

Row-level triggers can return one of the following:

`return;`

Execute the operation

Don't execute the operation

Indicates that the NEW row was modified by the trigger function.

```
CREATE TABLE test (
    i int,
    v varchar
);

CREATE OR REPLACE FUNCTION valid_id() RETURNS trigger AS $$
    if (($_TD->{new}{i} >= 100) || ($_TD->{new}{i} <= 0)) {
        return "SKIP";      # skip INSERT/UPDATE command
    } elsif ($_TD->{new}{v} ne "immortal") {
        $_TD->{new}{v} .= "(modified by trigger)";
        return "MODIFY";    # modify row and execute INSERT/UPDATE command
    } else {
        return;              # execute INSERT/UPDATE command
    }
$$ LANGUAGE plperl;

CREATE TRIGGER test_valid_id_trig
    BEFORE INSERT OR UPDATE ON test
    FOR EACH ROW EXECUTE FUNCTION valid_id();
```

PL/Perl can be used to write event trigger functions. In an event trigger function, the hash reference `$_TD` contains information about the current trigger event. `$_TD` is a global variable, which gets a separate local value for each invocation of the trigger. The fields of the `$_TD` hash reference are:

The name of the event the trigger is fired for.

The command tag for which the trigger is fired.

The return value of the trigger function is ignored.

```
CREATE OR REPLACE FUNCTION perlsnitch() RETURNS event_trigger AS $$
    elog(NOTICE, "perlsnitch: " . $_TD->{event} . " " . $_TD->{tag} . " ");
$$ LANGUAGE plperl;

CREATE EVENT TRIGGER perl_a_snitch
    ON ddl_command_start
    EXECUTE FUNCTION perlsnitch();
```

46.8.1. Configuration

1362

`plperl.on_init (string)`

Specifies Perl code to be executed when a Perl interpreter is first initialized, before it is specialized for use by `plperl` or `plperlu`. The SPI functions are not available when this code is executed. If the code fails with an error it will abort the initialization of the interpreter and propagate out to the calling query, causing the current transaction or subtransaction to be aborted.

The Perl code is limited to a single string. Longer code can be placed into a module and loaded by the `on_init` string. Examples:

```
plperl.on_init = 'require "plperlinit.pl"'
plperl.on_init = 'use lib "/my/app"; use MyApp::PgInit;'
```

Any modules loaded by `plperl.on_init`, either directly or indirectly, will be available for use by `plperl`. This may create a security risk. To see what modules have been loaded you can use:

```
DO 'elog(WARNING, join ", ", sort keys %INC)' LANGUAGE plperl;
```

Initialization will happen in the postmaster if the `plperl` library is included in [shared_preload_libraries](#), in which case extra consideration should be given to the risk of destabilizing the postmaster. The principal reason for making use of this feature is that Perl modules loaded by `plperl.on_init` need be loaded only at postmaster start, and will be instantly available without loading overhead in individual database sessions. However, keep in mind that the overhead is avoided only for the first Perl interpreter used by a database session — either PL/PerlU, or PL/Perl for the first SQL role that calls a PL/Perl function. Any additional Perl interpreters created in a database session will have to execute `plperl.on_init` afresh. Also, on Windows there will be no savings whatsoever from pre-loading, since the Perl interpreter created in the postmaster process does not propagate to child processes.

This parameter can only be set in the `postgresql.conf` file or on the server command line.

`plperl.on_plperl_init (string)`
`plperl.on_plperlu_init (string)`

These parameters specify Perl code to be executed when a Perl interpreter is specialized for `plperl` or `plperlu` respectively. This will happen when a PL/Perl or PL/PerlU function is first executed in a database session, or when an additional interpreter has to be created because the other language is called or a PL/Perl function is called by a new SQL role. This follows any initialization done by `plperl.on_init`. The SPI functions are not available when this code is executed. The Perl code in `plperl.on_plperl_init` is executed after “locking down” the interpreter, and thus it can only perform trusted operations.

If the code fails with an error it will abort the initialization and propagate out to the calling query, causing the current transaction or subtransaction to be aborted. Any actions already done within Perl won't be undone; however, that interpreter won't be used again. If the language is used again the initialization will be attempted again within a fresh Perl interpreter.

Only superusers can change these settings. Although these settings can be changed within a session, such changes will not affect Perl interpreters that have already been used to execute functions.

`plperl.use_strict (boolean)`

When set true subsequent compilations of PL/Perl functions will have the `strict` pragma enabled. This parameter does not affect functions already compiled in the current session.

46.8.2. Limitations and Missing Features

The following features are currently missing from PL/Perl, but they would make welcome contributions.

- PL/Perl functions cannot call each other directly.
- SPI is not yet fully implemented.

- If you are fetching very large data sets using `spi_exec_query`, you should be aware that these will all go into memory. You can avoid this by using `spi_query/spi_fetchrow` as illustrated earlier.

A similar problem occurs if a set-returning function passes a large set of rows back to Postgres Pro via `return`. You can avoid this problem too by instead using `return_next` for each row returned, as shown previously.

- When a session ends normally, not due to a fatal error, any `END` blocks that have been defined are executed. Currently no other actions are performed. Specifically, file handles are not automatically flushed and objects are not automatically destroyed.

Chapter 47. PL/Python — Python Procedural Language

The PL/Python procedural language allows Postgres Pro functions and procedures to be written in the *Python language*.

To install PL/Python in a particular database, use `CREATE EXTENSION plpython3u`.

Tip

If a language is installed into `template1`, all subsequently created databases will have the language installed automatically.

PL/Python is only available as an “untrusted” language, meaning it does not offer any way of restricting what users can do in it and is therefore named `plpython3u`. A trusted variant `plpython` might become available in the future if a secure execution mechanism is developed in Python. The writer of a function in untrusted PL/Python must take care that the function cannot be used to do anything unwanted, since it will be able to do anything that could be done by a user logged in as the database administrator. Only superusers can create functions in untrusted languages such as `plpython3u`.

Note

Users of source packages must specially enable the build of PL/Python during the installation process. (Refer to the installation instructions for more information.) Users of binary packages might find PL/Python in a separate subpackage.

47.1. PL/Python Functions

Functions in PL/Python are declared via the standard `CREATE FUNCTION` syntax:

```
CREATE FUNCTION funcname (argument-list)
    RETURNS return-type
AS $$
    # PL/Python function body
$$ LANGUAGE plpython3u;
```

The body of a function is simply a Python script. When the function is called, its arguments are passed as elements of the list `args`; named arguments are also passed as ordinary variables to the Python script. Use of named arguments is usually more readable. The result is returned from the Python code in the usual way, with `return` or `yield` (in case of a result-set statement). If you do not provide a return value, Python returns the default `None`. PL/Python translates Python's `None` into the SQL null value. In a procedure, the result from the Python code must be `None` (typically achieved by ending the procedure without a `return` statement or by using a `return` statement without argument); otherwise, an error will be raised.

For example, a function to return the greater of two integers can be defined as:

```
CREATE FUNCTION pymax (a integer, b integer)
    RETURNS integer
AS $$
    if a > b:
        return a
    return b
$$ LANGUAGE plpython3u;
```

The Python code that is given as the body of the function definition is transformed into a Python function. For example, the above results in:

```
def __plpython_procedure_pymax_23456():
    if a > b:
        return a
    return b
```

assuming that 23456 is the OID assigned to the function by Postgres Pro.

The arguments are set as global variables. Because of the scoping rules of Python, this has the subtle consequence that an argument variable cannot be reassigned inside the function to the value of an expression that involves the variable name itself, unless the variable is redeclared as global in the block. For example, the following won't work:

```
CREATE FUNCTION pystrip(x text)
    RETURNS text
AS $$
    x = x.strip() # error
    return x
$$ LANGUAGE plpython3u;
```

because assigning to `x` makes `x` a local variable for the entire block, and so the `x` on the right-hand side of the assignment refers to a not-yet-assigned local variable `x`, not the PL/Python function parameter. Using the `global` statement, this can be made to work:

```
CREATE FUNCTION pystrip(x text)
    RETURNS text
AS $$
    global x
    x = x.strip() # ok now
    return x
$$ LANGUAGE plpython3u;
```

But it is advisable not to rely on this implementation detail of PL/Python. It is better to treat the function parameters as read-only.

47.2. Data Values

Generally speaking, the aim of PL/Python is to provide a “natural” mapping between the PostgreSQL and the Python worlds. This informs the data mapping rules described below.

47.2.1. Data Type Mapping

When a PL/Python function is called, its arguments are converted from their PostgreSQL data type to a corresponding Python type:

- PostgreSQL `boolean` is converted to Python `bool`.
- PostgreSQL `smallint`, `int`, `bigint` and `oid` are converted to Python `int`.
- PostgreSQL `real` and `double` are converted to Python `float`.
- PostgreSQL `numeric` is converted to Python `Decimal`. This type is imported from the `cdecimal` package if that is available. Otherwise, `decimal.Decimal` from the standard library will be used. `cdecimal` is significantly faster than `decimal`. In Python 3.3 and up, however, `cdecimal` has been integrated into the standard library under the name `decimal`, so there is no longer any difference.
- PostgreSQL `bytea` is converted to Python `bytes`.
- All other data types, including the PostgreSQL character string types, are converted to a Python `str` (in Unicode like all Python strings).
- For nonscalar data types, see below.

When a PL/Python function returns, its return value is converted to the function's declared PostgreSQL return data type as follows:

- When the PostgreSQL return type is `boolean`, the return value will be evaluated for truth according to the *Python* rules. That is, 0 and empty string are false, but notably 'f' is true.
- When the PostgreSQL return type is `bytea`, the return value will be converted to Python `bytes` using the respective Python built-ins, with the result being converted to `bytea`.
- For all other PostgreSQL return types, the return value is converted to a string using the Python built-in `str`, and the result is passed to the input function of the PostgreSQL data type. (If the Python value is a `float`, it is converted using the `repr` built-in instead of `str`, to avoid loss of precision.)

Strings are automatically converted to the PostgreSQL server encoding when they are passed to PostgreSQL.

- For nonscalar data types, see below.

Note that logical mismatches between the declared PostgreSQL return type and the Python data type of the actual return object are not flagged; the value will be converted in any case.

47.2.2. Null, None

If an SQL null value is passed to a function, the argument value will appear as `None` in Python. For example, the function definition of `pymax` shown in [Section 47.1](#) will return the wrong answer for null inputs. We could add `STRICT` to the function definition to make Postgres Pro do something more reasonable: if a null value is passed, the function will not be called at all, but will just return a null result automatically. Alternatively, we could check for null inputs in the function body:

```
CREATE FUNCTION pyrax (a integer, b integer)
  RETURNS integer
AS $$
  if (a is None) or (b is None):
    return None
  if a > b:
    return a
  return b
$$ LANGUAGE plpython3u;
```

As shown above, to return an SQL null value from a PL/Python function, return the value `None`. This can be done whether the function is strict or not.

47.2.3. Arrays, Lists

SQL array values are passed into PL/Python as a Python list. To return an SQL array value out of a PL/Python function, return a Python list:

```
CREATE FUNCTION return_arr()
  RETURNS int[]
AS $$
  return [1, 2, 3, 4, 5]
$$ LANGUAGE plpython3u;
```

```
SELECT return_arr();
   return_arr
-----
 {1,2,3,4,5}
(1 row)
```

Multidimensional arrays are passed into PL/Python as nested Python lists. A 2-dimensional array is a list of lists, for example. When returning a multi-dimensional SQL array out of a PL/Python function, the inner lists at each level must all be of the same size. For example:

```
CREATE FUNCTION test_type_conversion_array_int4(x int4[]) RETURNS int4[] AS $$
```

```
plpy.info(x, type(x))
return x
$$ LANGUAGE plpython3u;

SELECT * FROM test_type_conversion_array_int4(ARRAY[[1,2,3],[4,5,6]]);
INFO:  ([[1, 2, 3], [4, 5, 6]], <type 'list'>)
test_type_conversion_array_int4
-----
{{1,2,3},{4,5,6}}
(1 row)
```

Other Python sequences, like tuples, are also accepted for backwards-compatibility with PostgreSQL versions 9.6 and below, when multi-dimensional arrays were not supported. However, they are always treated as one-dimensional arrays, because they are ambiguous with composite types. For the same reason, when a composite type is used in a multi-dimensional array, it must be represented by a tuple, rather than a list.

Note that in Python, strings are sequences, which can have undesirable effects that might be familiar to Python programmers:

```
CREATE FUNCTION return_str_arr()
  RETURNS varchar[]
AS $$
return "hello"
$$ LANGUAGE plpython3u;

SELECT return_str_arr();
return_str_arr
-----
{h,e,l,l,o}
(1 row)
```

47.2.4. Composite Types

Composite-type arguments are passed to the function as Python mappings. The element names of the mapping are the attribute names of the composite type. If an attribute in the passed row has the null value, it has the value `None` in the mapping. Here is an example:

```
CREATE TABLE employee (
  name text,
  salary integer,
  age integer
);

CREATE FUNCTION overpaid (e employee)
  RETURNS boolean
AS $$
if e["salary"] > 200000:
  return True
if (e["age"] < 30) and (e["salary"] > 100000):
  return True
return False
$$ LANGUAGE plpython3u;
```

There are multiple ways to return row or composite types from a Python function. The following examples assume we have:

```
CREATE TYPE named_value AS (
  name text,
  value integer
);
```


A composite result can be returned as a:

Sequence type (a tuple or list, but not a set because it is not indexable)

Returned sequence objects must have the same number of items as the composite result type has fields. The item with index 0 is assigned to the first field of the composite type, 1 to the second and so on. For example:

```
CREATE FUNCTION make_pair (name text, value integer)
  RETURNS named_value
AS $$
  return ( name, value )
  # or alternatively, as list: return [ name, value ]
$$ LANGUAGE plpython3u;
```

To return an SQL null for any column, insert `None` at the corresponding position.

When an array of composite types is returned, it cannot be returned as a list, because it is ambiguous whether the Python list represents a composite type, or another array dimension.

Mapping (dictionary)

The value for each result type column is retrieved from the mapping with the column name as key. Example:

```
CREATE FUNCTION make_pair (name text, value integer)
  RETURNS named_value
AS $$
  return { "name": name, "value": value }
$$ LANGUAGE plpython3u;
```

Any extra dictionary key/value pairs are ignored. Missing keys are treated as errors. To return an SQL null value for any column, insert `None` with the corresponding column name as the key.

Object (any object providing method `__getattr__`)

This works the same as a mapping. Example:

```
CREATE FUNCTION make_pair (name text, value integer)
  RETURNS named_value
AS $$
  class named_value:
    def __init__ (self, n, v):
      self.name = n
      self.value = v
  return named_value(name, value)

  # or simply
  class nv: pass
  nv.name = name
  nv.value = value
  return nv
$$ LANGUAGE plpython3u;
```

Functions with `OUT` parameters are also supported. For example:

```
CREATE FUNCTION multiout_simple(OUT i integer, OUT j integer) AS $$
  return (1, 2)
$$ LANGUAGE plpython3u;
```

```
SELECT * FROM multiout_simple();
```

Output parameters of procedures are passed back the same way. For example:

```
CREATE PROCEDURE python_triple(INOUT a integer, INOUT b integer) AS $$
```

```
return (a * 3, b * 3)
$$ LANGUAGE plpython3u;

CALL python_triple(5, 10);
```

47.2.5. Set-Returning Functions

A PL/Python function can also return sets of scalar or composite types. There are several ways to achieve this because the returned object is internally turned into an iterator. The following examples assume we have composite type:

```
CREATE TYPE greeting AS (
    how text,
    who text
);
```

A set result can be returned from a:

Sequence type (tuple, list, set)

```
CREATE FUNCTION greet (how text)
    RETURNS SETOF greeting
AS $$
    # return tuple containing lists as composite types
    # all other combinations work also
    return ( [ how, "World" ], [ how, "PostgreSQL" ], [ how, "PL/Python" ] )
$$ LANGUAGE plpython3u;
```

Iterator (any object providing `__iter__` and `__next__` methods)

```
CREATE FUNCTION greet (how text)
    RETURNS SETOF greeting
AS $$
    class producer:
        def __init__ (self, how, who):
            self.how = how
            self.who = who
            self.ndx = -1

        def __iter__ (self):
            return self

        def __next__(self):
            self.ndx += 1
            if self.ndx == len(self.who):
                raise StopIteration
            return ( self.how, self.who[self.ndx] )

    return producer(how, [ "World", "PostgreSQL", "PL/Python" ])
$$ LANGUAGE plpython3u;
```

Generator (yield)

```
CREATE FUNCTION greet (how text)
    RETURNS SETOF greeting
AS $$
    for who in [ "World", "PostgreSQL", "PL/Python" ]:
        yield ( how, who )
$$ LANGUAGE plpython3u;
```

Set-returning functions with `OUT` parameters (using `RETURNS SETOF record`) are also supported. For example:

```
CREATE FUNCTION multiout_simple_setof(n integer, OUT integer, OUT integer) RETURNS
    SETOF record AS $$
return [(1, 2)] * n
$$ LANGUAGE plpython3u;

SELECT * FROM multiout_simple_setof(3);
```

47.3. Sharing Data

The global dictionary `SD` is available to store private data between repeated calls to the same function. The global dictionary `GD` is public data, that is available to all Python functions within a session; use with care.

Each function gets its own execution environment in the Python interpreter, so that global data and function arguments from `myfunc` are not available to `myfunc2`. The exception is the data in the `GD` dictionary, as mentioned above.

47.4. Anonymous Code Blocks

PL/Python also supports anonymous code blocks called with the `DO` statement:

```
DO $$
    # PL/Python code
$$ LANGUAGE plpython3u;
```

An anonymous code block receives no arguments, and whatever value it might return is discarded. Otherwise it behaves just like a function.

47.5. Trigger Functions

When a function is used as a trigger, the dictionary `TD` contains trigger-related values:

```
TD["event"]
    contains the event as a string: INSERT, UPDATE, DELETE, or TRUNCATE.
```

```
TD["when"]
    contains one of BEFORE, AFTER, or INSTEAD OF.
```

```
TD["level"]
    contains ROW or STATEMENT.
```

```
TD["new"]
TD["old"]
    For a row-level trigger, one or both of these fields contain the respective trigger rows, depending on the trigger event.
```

```
TD["name"]
    contains the trigger name.
```

```
TD["table_name"]
    contains the name of the table on which the trigger occurred.
```

```
TD["table_schema"]
    contains the schema of the table on which the trigger occurred.
```

```
TD["relid"]
    contains the OID of the table on which the trigger occurred.
```

`TD["args"]`

If the `CREATE TRIGGER` command included arguments, they are available in `TD["args"][0]` to `TD["args"][n-1]`.

If `TD["when"]` is `BEFORE` or `INSTEAD OF` and `TD["level"]` is `ROW`, you can return `None` or `"OK"` from the Python function to indicate the row is unmodified, `"SKIP"` to abort the event, or if `TD["event"]` is `INSERT` or `UPDATE` you can return `"MODIFY"` to indicate you've modified the new row. Otherwise the return value is ignored.

47.6. Database Access

The PL/Python language module automatically imports a Python module called `plpy`. The functions and constants in this module are available to you in the Python code as `plpy.foo`.

47.6.1. Database Access Functions

The `plpy` module provides several functions to execute database commands:

```
plpy.execute(query [, limit])
```

Calling `plpy.execute` with a query string and an optional row limit argument causes that query to be run and the result to be returned in a result object.

If `limit` is specified and is greater than zero, then `plpy.execute` retrieves at most `limit` rows, much as if the query included a `LIMIT` clause. Omitting `limit` or specifying it as zero results in no row limit.

The result object emulates a list or dictionary object. The result object can be accessed by row number and column name. For example:

```
rv = plpy.execute("SELECT * FROM my_table", 5)
```

returns up to 5 rows from `my_table`. If `my_table` has a column `my_column`, it would be accessed as:

```
foo = rv[i]["my_column"]
```

The number of rows returned can be obtained using the built-in `len` function.

The result object has these additional methods:

```
nrows()
```

Returns the number of rows processed by the command. Note that this is not necessarily the same as the number of rows returned. For example, an `UPDATE` command will set this value but won't return any rows (unless `RETURNING` is used).

```
status()
```

The `SPI_execute()` return value.

```
colnames()
```

```
coltypes()
```

```
coltypmods()
```

Return a list of column names, list of column type OIDs, and list of type-specific type modifiers for the columns, respectively.

These methods raise an exception when called on a result object from a command that did not produce a result set, e.g., `UPDATE` without `RETURNING`, or `DROP TABLE`. But it is OK to use these methods on a result set containing zero rows.

```
__str__()
```

The standard `__str__` method is defined so that it is possible for example to debug query execution results using `plpy.debug(rv)`.

The result object can be modified.

Note that calling `plpy.execute` will cause the entire result set to be read into memory. Only use that function when you are sure that the result set will be relatively small. If you don't want to risk excessive memory usage when fetching large results, use `plpy.cursor` rather than `plpy.execute`.

```
plpy.prepare(query [, argtypes])
plpy.execute(plan [, arguments [, limit]])
```

`plpy.prepare` prepares the execution plan for a query. It is called with a query string and a list of parameter types, if you have parameter references in the query. For example:

```
plan = plpy.prepare("SELECT last_name FROM my_users WHERE first_name = $1",
    ["text"])
```

`text` is the type of the variable you will be passing for `$1`. The second argument is optional if you don't want to pass any parameters to the query.

After preparing a statement, you use a variant of the function `plpy.execute` to run it:

```
rv = plpy.execute(plan, ["name"], 5)
```

Pass the plan as the first argument (instead of the query string), and a list of values to substitute into the query as the second argument. The second argument is optional if the query does not expect any parameters. The third argument is the optional row limit as before.

Alternatively, you can call the `execute` method on the plan object:

```
rv = plan.execute(["name"], 5)
```

Query parameters and result row fields are converted between PostgreSQL and Python data types as described in [Section 47.2](#).

When you prepare a plan using the PL/Python module it is automatically saved. Read the SPI documentation ([Chapter 48](#)) for a description of what this means. In order to make effective use of this across function calls one needs to use one of the persistent storage dictionaries `SD` or `GD` (see [Section 47.3](#)). For example:

```
CREATE FUNCTION usesavedplan() RETURNS trigger AS $$
    if "plan" in SD:
        plan = SD["plan"]
    else:
        plan = plpy.prepare("SELECT 1")
        SD["plan"] = plan
    # rest of function
$$ LANGUAGE plpython3u;
```

```
plpy.cursor(query)
plpy.cursor(plan [, arguments])
```

The `plpy.cursor` function accepts the same arguments as `plpy.execute` (except for the row limit) and returns a cursor object, which allows you to process large result sets in smaller chunks. As with `plpy.execute`, either a query string or a plan object along with a list of arguments can be used, or the `cursor` function can be called as a method of the plan object.

The cursor object provides a `fetch` method that accepts an integer parameter and returns a result object. Each time you call `fetch`, the returned object will contain the next batch of rows, never larger than the parameter value. Once all rows are exhausted, `fetch` starts returning an empty result object. Cursor objects also provide an [iterator interface](#), yielding one row at a time until all rows are exhausted. Data fetched that way is not returned as result objects, but rather as dictionaries, each dictionary corresponding to a single result row.

An example of two ways of processing data from a large table is:

```
CREATE FUNCTION count_odd_iterator() RETURNS integer AS $$
odd = 0
for row in plpy.cursor("select num from targetable"):
    if row['num'] % 2:
        odd += 1
return odd
$$ LANGUAGE plpython3u;

CREATE FUNCTION count_odd_fetch(batch_size integer) RETURNS integer AS $$
odd = 0
cursor = plpy.cursor("select num from targetable")
while True:
    rows = cursor.fetch(batch_size)
    if not rows:
        break
    for row in rows:
        if row['num'] % 2:
            odd += 1
return odd
$$ LANGUAGE plpython3u;

CREATE FUNCTION count_odd_prepared() RETURNS integer AS $$
odd = 0
plan = plpy.prepare("select num from targetable where num % $1 <> 0", ["integer"])
rows = list(plpy.cursor(plan, [2])) # or: = list(plan.cursor([2]))

return len(rows)
$$ LANGUAGE plpython3u;
```

Cursors are automatically disposed of. But if you want to explicitly release all resources held by a cursor, use the `close` method. Once closed, a cursor cannot be fetched from anymore.

Tip

Do not confuse objects created by `plpy.cursor` with DB-API cursors as defined by the [Python Database API specification](#). They don't have anything in common except for the name.

47.6.2. Trapping Errors

Functions accessing the database might encounter errors, which will cause them to abort and raise an exception. Both `plpy.execute` and `plpy.prepare` can raise an instance of a subclass of `plpy.SPIError`, which by default will terminate the function. This error can be handled just like any other Python exception, by using the `try/except` construct. For example:

```
CREATE FUNCTION try_adding_joe() RETURNS text AS $$
try:
    plpy.execute("INSERT INTO users(username) VALUES ('joe')")
except plpy.SPIError:
    return "something went wrong"
else:
    return "Joe added"
$$ LANGUAGE plpython3u;
```

The actual class of the exception being raised corresponds to the specific condition that caused the error. Refer to [Table A.1](#) for a list of possible conditions. The module `plpy.spiexceptions` defines an exception class for each Postgres Pro condition, deriving their names from the condition name. For instance, `division_by_zero` becomes `DivisionByZero`, `unique_violation` becomes `UniqueViolation`,

`fdw_error` becomes `FdwError`, and so on. Each of these exception classes inherits from `SPLError`. This separation makes it easier to handle specific errors, for instance:

```
CREATE FUNCTION insert_fraction(numerator int, denominator int) RETURNS text AS $$
from plpy import spiexceptions
try:
    plan = plpy.prepare("INSERT INTO fractions (frac) VALUES ($1 / $2)", ["int",
    "int"])
    plpy.execute(plan, [numerator, denominator])
except spiexceptions.DivisionByZero:
    return "denominator cannot equal zero"
except spiexceptions.UniqueViolation:
    return "already have that fraction"
except plpy.SPLError as e:
    return "other error, SQLSTATE %s" % e.sqlstate
else:
    return "fraction inserted"
$$ LANGUAGE plpython3u;
```

Note that because all exceptions from the `plpy.spiexceptions` module inherit from `SPLError`, an `except` clause handling it will catch any database access error.

As an alternative way of handling different error conditions, you can catch the `SPLError` exception and determine the specific error condition inside the `except` block by looking at the `sqlstate` attribute of the exception object. This attribute is a string value containing the “SQLSTATE” error code. This approach provides approximately the same functionality

47.7. Explicit Subtransactions

Recovering from errors caused by database access as described in [Section 47.6.2](#) can lead to an undesirable situation where some operations succeed before one of them fails, and after recovering from that error the data is left in an inconsistent state. PL/Python offers a solution to this problem in the form of explicit subtransactions.

47.7.1. Subtransaction Context Managers

Consider a function that implements a transfer between two accounts:

```
CREATE FUNCTION transfer_funds() RETURNS void AS $$
try:
    plpy.execute("UPDATE accounts SET balance = balance - 100 WHERE account_name =
    'joe'")
    plpy.execute("UPDATE accounts SET balance = balance + 100 WHERE account_name =
    'mary'")
except plpy.SPLError as e:
    result = "error transferring funds: %s" % e.args
else:
    result = "funds transferred correctly"
plan = plpy.prepare("INSERT INTO operations (result) VALUES ($1)", ["text"])
plpy.execute(plan, [result])
$$ LANGUAGE plpython3u;
```

If the second `UPDATE` statement results in an exception being raised, this function will report the error, but the result of the first `UPDATE` will nevertheless be committed. In other words, the funds will be withdrawn from Joe's account, but will not be transferred to Mary's account.

To avoid such issues, you can wrap your `plpy.execute` calls in an explicit subtransaction. The `plpy` module provides a helper object to manage explicit subtransactions that gets created with the `plpy.subtransaction()` function. Objects created by this function implement the [context manager interface](#). Using explicit subtransactions we can rewrite our function as:

```
CREATE FUNCTION transfer_funds2() RETURNS void AS $$
try:
    with plpy.subtransaction():
        plpy.execute("UPDATE accounts SET balance = balance - 100 WHERE account_name =
'joe'")
        plpy.execute("UPDATE accounts SET balance = balance + 100 WHERE account_name =
'mary'")
except plpy.SPIError as e:
    result = "error transferring funds: %s" % e.args
else:
    result = "funds transferred correctly"
plan = plpy.prepare("INSERT INTO operations (result) VALUES ($1)", ["text"])
plpy.execute(plan, [result])
$$ LANGUAGE plpython3u;
```

Note that the use of `try/except` is still required. Otherwise the exception would propagate to the top of the Python stack and would cause the whole function to abort with a Postgres Pro error, so that the `operations` table would not have any row inserted into it. The subtransaction context manager does not trap errors, it only assures that all database operations executed inside its scope will be atomically committed or rolled back. A rollback of the subtransaction block occurs on any kind of exception exit, not only ones caused by errors originating from database access. A regular Python exception raised inside an explicit subtransaction block would also cause the subtransaction to be rolled back.

47.8. Transaction Management

In a procedure called from the top level or an anonymous code block (`DO` command) called from the top level it is possible to control transactions. To commit the current transaction, call `plpy.commit()`. To roll back the current transaction, call `plpy.rollback()`. (Note that it is not possible to run the SQL commands `COMMIT` or `ROLLBACK` via `plpy.execute` or similar. It has to be done using these functions.) After a transaction is ended, a new transaction is automatically started, so there is no separate function for that.

Here is an example:

```
CREATE PROCEDURE transaction_test1()
LANGUAGE plpython3u
AS $$
for i in range(0, 10):
    plpy.execute("INSERT INTO test1 (a) VALUES (%d)" % i)
    if i % 2 == 0:
        plpy.commit()
    else:
        plpy.rollback()
$$;

CALL transaction_test1();
```

Transactions cannot be ended when an explicit subtransaction is active.

47.9. Utility Functions

The `plpy` module also provides the functions

```
plpy.debug( msg, **kwargs )
plpy.log( msg, **kwargs )
plpy.info( msg, **kwargs )
plpy.notice( msg, **kwargs )
plpy.warning( msg, **kwargs )
plpy.error( msg, **kwargs )
plpy.fatal( msg, **kwargs )
```


`plpy.error` and `plpy.fatal` actually raise a Python exception which, if uncaught, propagates out to the calling query, causing the current transaction or subtransaction to be aborted. `raise plpy.Error(msg)` and `raise plpy.Fatal(msg)` are equivalent to calling `plpy.error(msg)` and `plpy.fatal(msg)`, respectively but the `raise` form does not allow passing keyword arguments. The other functions only generate messages of different priority levels. Whether messages of a particular priority are reported to the client, written to the server log, or both is controlled by the [log_min_messages](#) and [client_min_messages](#) configuration variables. See [Chapter 19](#) for more information.

The `msg` argument is given as a positional argument. For backward compatibility, more than one positional argument can be given. In that case, the string representation of the tuple of positional arguments becomes the message reported to the client.

The following keyword-only arguments are accepted:

```
detail
hint
sqlstate
schema_name
table_name
column_name
datatype_name
constraint_name
```

The string representation of the objects passed as keyword-only arguments is used to enrich the messages reported to the client. For example:

```
CREATE FUNCTION raise_custom_exception() RETURNS void AS $$
plpy.error("custom exception message",
           detail="some info about exception",
           hint="hint for users")
$$ LANGUAGE plpython3u;

=# SELECT raise_custom_exception();
ERROR:  plpy.Error: custom exception message
DETAIL:  some info about exception
HINT:   hint for users
CONTEXT:  Traceback (most recent call last):
  PL/Python function "raise_custom_exception", line 4, in <module>
    hint="hint for users")
PL/Python function "raise_custom_exception"
```

Another set of utility functions are `plpy.quote_literal(string)`, `plpy.quote_nullable(string)`, and `plpy.quote_ident(string)`. They are equivalent to the built-in quoting functions described in [Section 9.4](#). They are useful when constructing ad-hoc queries. A PL/Python equivalent of dynamic SQL from [Example 44.1](#) would be:

```
plpy.execute("UPDATE tbl SET %s = %s WHERE key = %s" % (
    plpy.quote_ident(colname),
    plpy.quote_nullable(newvalue),
    plpy.quote_literal(keyvalue)))
```

47.10. Python 2 vs. Python 3

PL/Python supports only Python 3. Past versions of Postgres Pro supported Python 2, using the `plpythonu` and `plpython2u` language names.

47.11. Environment Variables

Some of the environment variables that are accepted by the Python interpreter can also be used to affect PL/Python behavior. They would need to be set in the environment of the main Postgres Pro server

process, for example in a start script. The available environment variables depend on the version of Python; see the Python documentation for details. At the time of this writing, the following environment variables have an affect on PL/Python, assuming an adequate Python version:

- PYTHONHOME
- PYTHONPATH
- PYTHON2K
- PYTHONOPTIMIZE
- PYTHONDEBUG
- PYTHONVERBOSE
- PYTHONCASEOK
- PYTHONDONTWRITEBYTECODE
- PYTHONIOENCODING
- PYTHONUSERBASE
- PYTHONHASHSEED

(It appears to be a Python implementation detail beyond the control of PL/Python that some of the environment variables listed on the `python` man page are only effective in a command-line interpreter and not an embedded Python interpreter.)

Chapter 48. Server Programming Interface

The *Server Programming Interface* (SPI) gives writers of user-defined C functions the ability to run SQL commands inside their functions or procedures. SPI is a set of interface functions to simplify access to the parser, planner, and executor. SPI also does some memory management.

Note

The available procedural languages provide various means to execute SQL commands from functions. Most of these facilities are based on SPI, so this documentation might be of use for users of those languages as well.

Note that if a command invoked via SPI fails, then control will not be returned to your C function. Rather, the transaction or subtransaction in which your C function executes will be rolled back. (This might seem surprising given that the SPI functions mostly have documented error-return conventions. Those conventions only apply for errors detected within the SPI functions themselves, however.) It is possible to recover control after an error by establishing your own subtransaction surrounding SPI calls that might fail.

SPI functions return a nonnegative result on success (either via a returned integer value or in the global variable `SPI_result`, as described below). On error, a negative result or `NULL` will be returned.

Source code files that use SPI must include the header file `executor/spi.h`.

48.1. Interface Functions

SPI_connect

`SPI_connect`, `SPI_connect_ext` — connect a C function to the SPI manager

Synopsis

```
int SPI_connect(void)
int SPI_connect_ext(int options)
```

Description

`SPI_connect` opens a connection from a C function invocation to the SPI manager. You must call this function if you want to execute commands through SPI. Some utility SPI functions can be called from unconnected C functions.

`SPI_connect_ext` does the same but has an argument that allows passing option flags. Currently, the following option values are available:

`SPI_OPT_NONATOMIC`

Sets the SPI connection to be *nonatomic*, which means that transaction control calls (`SPI_commit`, `SPI_rollback`) are allowed. Otherwise, calling those functions will result in an immediate error.

`SPI_connect()` is equivalent to `SPI_connect_ext(0)`.

Return Value

`SPI_OK_CONNECT`

on success

`SPI_ERROR_CONNECT`

on error

SPI_finish

SPI_finish — disconnect a C function from the SPI manager

Synopsis

```
int SPI_finish(void)
```

Description

SPI_finish closes an existing connection to the SPI manager. You must call this function after completing the SPI operations needed during your C function's current invocation. You do not need to worry about making this happen, however, if you abort the transaction via `elog(ERROR)`. In that case SPI will clean itself up automatically.

Return Value

SPI_OK_FINISH

if properly disconnected

SPI_ERROR_UNCONNECTED

if called from an unconnected C function

SPI_execute

SPI_execute — execute a command

Synopsis

```
int SPI_execute(const char * command, bool read_only, long count)
```

Description

SPI_execute executes the specified SQL command for *count* rows. If *read_only* is true, the command must be read-only, and execution overhead is somewhat reduced.

This function can only be called from a connected C function.

If *count* is zero then the command is executed for all rows that it applies to. If *count* is greater than zero, then no more than *count* rows will be retrieved; execution stops when the count is reached, much like adding a LIMIT clause to the query. For example,

```
SPI_execute("SELECT * FROM foo", true, 5);
```

will retrieve at most 5 rows from the table. Note that such a limit is only effective when the command actually returns rows. For example,

```
SPI_execute("INSERT INTO foo SELECT * FROM bar", false, 5);
```

inserts all rows from bar, ignoring the *count* parameter. However, with

```
SPI_execute("INSERT INTO foo SELECT * FROM bar RETURNING *", false, 5);
```

at most 5 rows would be inserted, since execution would stop after the fifth RETURNING result row is retrieved.

You can pass multiple commands in one string; SPI_execute returns the result for the command executed last. The *count* limit applies to each command separately (even though only the last result will actually be returned). The limit is not applied to any hidden commands generated by rules.

When *read_only* is false, SPI_execute increments the command counter and computes a new *snapshot* before executing each command in the string. The snapshot does not actually change if the current transaction isolation level is SERIALIZABLE or REPEATABLE READ, but in READ COMMITTED mode the snapshot update allows each command to see the results of newly committed transactions from other sessions. This is essential for consistent behavior when the commands are modifying the database.

When *read_only* is true, SPI_execute does not update either the snapshot or the command counter, and it allows only plain SELECT commands to appear in the command string. The commands are executed using the snapshot previously established for the surrounding query. This execution mode is somewhat faster than the read/write mode due to eliminating per-command overhead. It also allows genuinely *stable* functions to be built: since successive executions will all use the same snapshot, there will be no change in the results.

It is generally unwise to mix read-only and read-write commands within a single function using SPI; that could result in very confusing behavior, since the read-only queries would not see the results of any database updates done by the read-write queries.

The actual number of rows for which the (last) command was executed is returned in the global variable SPI_processed. If the return value of the function is SPI_OK_SELECT, SPI_OK_INSERT_RETURNING, SPI_OK_DELETE_RETURNING, SPI_OK_UPDATE_RETURNING, or SPI_OK_MERGE_RETURNING, then you can use the global pointer SPITupleTable *SPI_tuptable to access the result rows. Some utility commands (such as EXPLAIN) also return row sets, and SPI_tuptable will contain the result in these cases too. Some utility commands (COPY, CREATE TABLE AS) don't return a row set, so SPI_tuptable is NULL, but they still return the number of rows processed in SPI_processed.

The structure `SPITupleTable` is defined thus:

```
typedef struct SPITupleTable
{
    /* Public members */
    TupleDesc    tupdesc;          /* tuple descriptor */
    HeapTuple    *vals;            /* array of tuples */
    uint64        numvals;         /* number of valid tuples */

    /* Private members, not intended for external callers */
    uint64        allocated;       /* allocated length of vals array */
    MemoryContext tupabcxt;        /* memory context of result table */
    slist_node    next;           /* link for internal bookkeeping */
    SubTransactionId subid;       /* subxact in which tuptable was created */
} SPITupleTable;
```

The fields `tupdesc`, `vals`, and `numvals` can be used by SPI callers; the remaining fields are internal. `vals` is an array of pointers to rows. The number of rows is given by `numvals` (for somewhat historical reasons, this count is also returned in `SPI_processed`). `tupdesc` is a row descriptor which you can pass to SPI functions dealing with rows.

`SPI_finish` frees all `SPITupleTables` allocated during the current C function. You can free a particular result table earlier, if you are done with it, by calling `SPI_freetupable`.

Arguments

```
const char * command
    string containing command to execute

bool read_only
    true for read-only execution

long count
    maximum number of rows to return, or 0 for no limit
```

Return Value

If the execution of the command was successful then one of the following (nonnegative) values will be returned:

```
SPI_OK_SELECT
    if a SELECT (but not SELECT INTO) was executed

SPI_OK_SELINTO
    if a SELECT INTO was executed

SPI_OK_INSERT
    if an INSERT was executed

SPI_OK_DELETE
    if a DELETE was executed

SPI_OK_UPDATE
    if an UPDATE was executed

SPI_OK_MERGE
    if a MERGE was executed
```

`SPI_OK_INSERT_RETURNING`

if an `INSERT RETURNING` was executed

`SPI_OK_DELETE_RETURNING`

if a `DELETE RETURNING` was executed

`SPI_OK_UPDATE_RETURNING`

if an `UPDATE RETURNING` was executed

`SPI_OK_MERGE_RETURNING`

if a `MERGE RETURNING` was executed

`SPI_OK_UTILITY`

if a utility command (e.g., `CREATE TABLE`) was executed

`SPI_OK_REWRITTEN`

if the command was rewritten into another kind of command (e.g., `UPDATE` became an `INSERT`) by a [rule](#).

On error, one of the following negative values is returned:

`SPI_ERROR_ARGUMENT`

if *command* is `NULL` or *count* is less than 0

`SPI_ERROR_COPY`

if `COPY TO stdout` or `COPY FROM stdin` was attempted

`SPI_ERROR_TRANSACTION`

if a transaction manipulation command was attempted (`BEGIN`, `COMMIT`, `ROLLBACK`, `SAVEPOINT`, `PREPARE TRANSACTION`, `COMMIT PREPARED`, `ROLLBACK PREPARED`, or any variant thereof)

`SPI_ERROR_OPUNKNOWN`

if the command type is unknown (shouldn't happen)

`SPI_ERROR_UNCONNECTED`

if called from an unconnected C function

Notes

All SPI query-execution functions set both `SPI_processed` and `SPI_tuptable` (just the pointer, not the contents of the structure). Save these two global variables into local C function variables if you need to access the result table of `SPI_execute` or another query-execution function across later calls.

SPI_exec

SPI_exec — execute a read/write command

Synopsis

```
int SPI_exec(const char * command, long count)
```

Description

SPI_exec is the same as SPI_execute, with the latter's *read_only* parameter always taken as false.

Arguments

`const char * command`

string containing command to execute

`long count`

maximum number of rows to return, or 0 for no limit

Return Value

See SPI_execute.

SPI_execute_extended

SPI_execute_extended — execute a command with out-of-line parameters

Synopsis

```
int SPI_execute_extended(const char *command,
                        const SPIExecuteOptions * options)
```

Description

SPI_execute_extended executes a command that might include references to externally supplied parameters. The command text refers to a parameter as *\$n*, and the *options->params* object (if supplied) provides values and type information for each such symbol. Various execution options can be specified in the *options* struct, too.

The *options->params* object should normally mark each parameter with the `PARAM_FLAG_CONST` flag, since a one-shot plan is always used for the query.

If *options->dest* is not NULL, then result tuples are passed to that object as they are generated by the executor, instead of being accumulated in `SPI_tuptable`. Using a caller-supplied `DestReceiver` object is particularly helpful for queries that might generate many tuples, since the data can be processed on-the-fly instead of being accumulated in memory.

Arguments

`const char * command`

command string

`const SPIExecuteOptions * options`

struct containing optional arguments

Callers should always zero out the entire *options* struct, then fill whichever fields they want to set. This ensures forward compatibility of code, since any fields that are added to the struct in future will be defined to behave backwards-compatibly if they are zero. The currently available *options* fields are:

`ParamListInfo params`

data structure containing query parameter types and values; NULL if none

`bool read_only`

true for read-only execution

`bool allow_nonatomic`

true allows non-atomic execution of CALL and DO statements (but this field is ignored unless the `SPI_OPT_NONATOMIC` flag was passed to `SPI_connect_ext`)

`bool must_return_tuples`

if true, raise error if the query is not of a kind that returns tuples (this does not forbid the case where it happens to return zero tuples)

`uint64 tcount`

maximum number of rows to return, or 0 for no limit

`DestReceiver * dest`

`DestReceiver` object that will receive any tuples emitted by the query; if NULL, result tuples are accumulated into a `SPI_tuptable` structure, as in `SPI_execute`

ResourceOwner *owner*

This field is present for consistency with `SPI_execute_plan_extended`, but it is ignored, since the plan used by `SPI_execute_extended` is never saved.

Return Value

The return value is the same as for `SPI_execute`.

When `options->dest` is `NULL`, `SPI_processed` and `SPI_tuptable` are set as in `SPI_execute`. When `options->dest` is not `NULL`, `SPI_processed` is set to zero and `SPI_tuptable` is set to `NULL`. If a tuple count is required, the caller's `DestReceiver` object must calculate it.

SPI_execute_with_args

`SPI_execute_with_args` — execute a command with out-of-line parameters

Synopsis

```
int SPI_execute_with_args(const char *command,
                          int nargs, Oid *argtypes,
                          Datum *values, const char *nulls,
                          bool read_only, long count)
```

Description

`SPI_execute_with_args` executes a command that might include references to externally supplied parameters. The command text refers to a parameter as `$n`, and the call specifies data types and values for each such symbol. `read_only` and `count` have the same interpretation as in `SPI_execute`.

The main advantage of this routine compared to `SPI_execute` is that data values can be inserted into the command without tedious quoting/escaping, and thus with much less risk of SQL-injection attacks.

Similar results can be achieved with `SPI_prepare` followed by `SPI_execute_plan`; however, when using this function the query plan is always customized to the specific parameter values provided. For one-time query execution, this function should be preferred. If the same command is to be executed with many different parameters, either method might be faster, depending on the cost of re-planning versus the benefit of custom plans.

Arguments

`const char * command`

command string

`int nargs`

number of input parameters (\$1, \$2, etc.)

`Oid * argtypes`

an array of length `nargs`, containing the OIDs of the data types of the parameters

`Datum * values`

an array of length `nargs`, containing the actual parameter values

`const char * nulls`

an array of length `nargs`, describing which parameters are null

If `nulls` is `NULL` then `SPI_execute_with_args` assumes that no parameters are null. Otherwise, each entry of the `nulls` array should be `' '` if the corresponding parameter value is non-null, or `'n'` if the corresponding parameter value is null. (In the latter case, the actual value in the corresponding `values` entry doesn't matter.) Note that `nulls` is not a text string, just an array: it does not need a `'\0'` terminator.

`bool read_only`

true for read-only execution

`long count`

maximum number of rows to return, or 0 for no limit

Return Value

The return value is the same as for `SPI_execute`.

`SPI_processed` and `SPI_tuptable` are set as in `SPI_execute` if successful.

SPI_prepare

SPI_prepare — prepare a statement, without executing it yet

Synopsis

```
SPIPlanPtr SPI_prepare(const char * command, int nargs, Oid * argtypes)
```

Description

SPI_prepare creates and returns a prepared statement for the specified command, but doesn't execute the command. The prepared statement can later be executed repeatedly using SPI_execute_plan.

When the same or a similar command is to be executed repeatedly, it is generally advantageous to perform parse analysis only once, and might furthermore be advantageous to re-use an execution plan for the command. SPI_prepare converts a command string into a prepared statement that encapsulates the results of parse analysis. The prepared statement also provides a place for caching an execution plan if it is found that generating a custom plan for each execution is not helpful.

A prepared command can be generalized by writing parameters (\$1, \$2, etc.) in place of what would be constants in a normal command. The actual values of the parameters are then specified when SPI_execute_plan is called. This allows the prepared command to be used over a wider range of situations than would be possible without parameters.

The statement returned by SPI_prepare can be used only in the current invocation of the C function, since SPI_finish frees memory allocated for such a statement. But the statement can be saved for longer using the functions SPI_keepplan or SPI_saveplan.

Arguments

const char * *command*

command string

int *nargs*

number of input parameters (\$1, \$2, etc.)

Oid * *argtypes*

pointer to an array containing the OIDs of the data types of the parameters

Return Value

SPI_prepare returns a non-null pointer to an SPIPlan, which is an opaque struct representing a prepared statement. On error, NULL will be returned, and SPI_result will be set to one of the same error codes used by SPI_execute, except that it is set to SPI_ERROR_ARGUMENT if *command* is NULL, or if *nargs* is less than 0, or if *nargs* is greater than 0 and *argtypes* is NULL.

Notes

If no parameters are defined, a generic plan will be created at the first use of SPI_execute_plan, and used for all subsequent executions as well. If there are parameters, the first few uses of SPI_execute_plan will generate custom plans that are specific to the supplied parameter values. After enough uses of the same prepared statement, SPI_execute_plan will build a generic plan, and if that is not too much more expensive than the custom plans, it will start using the generic plan instead of re-planning each time. If this default behavior is unsuitable, you can alter it by passing the CURSOR_OPT_GENERIC_PLAN or CURSOR_OPT_CUSTOM_PLAN flag to SPI_prepare_cursor, to force use of generic or custom plans respectively.

Although the main point of a prepared statement is to avoid repeated parse analysis and planning of the statement, Postgres Pro will force re-analysis and re-planning of the statement before using it whenever

database objects used in the statement have undergone definitional (DDL) changes since the previous use of the prepared statement. Also, if the value of `search_path` changes from one use to the next, the statement will be re-parsed using the new `search_path`. (This latter behavior is new as of PostgreSQL 9.3.) See [PREPARE](#) for more information about the behavior of prepared statements.

This function should only be called from a connected C function.

`SPIPlanPtr` is declared as a pointer to an opaque struct type in `spi.h`. It is unwise to try to access its contents directly, as that makes your code much more likely to break in future revisions of Postgres Pro.

The name `SPIPlanPtr` is somewhat historical, since the data structure no longer necessarily contains an execution plan.

SPI_prepare_cursor

SPI_prepare_cursor — prepare a statement, without executing it yet

Synopsis

```
SPIPlanPtr SPI_prepare_cursor(const char * command, int nargs,
                             Oid * argtypes, int cursorOptions)
```

Description

SPI_prepare_cursor is identical to SPI_prepare, except that it also allows specification of the planner's "cursor options" parameter. This is a bit mask having the values shown in `nodes/parsenodes.h` for the `options` field of `DeclareCursorStmt`. SPI_prepare always takes the cursor options as zero.

This function is now deprecated in favor of SPI_prepare_extended.

Arguments

`const char * command`
command string

`int nargs`
number of input parameters (\$1, \$2, etc.)

`Oid * argtypes`
pointer to an array containing the OIDs of the data types of the parameters

`int cursorOptions`
integer bit mask of cursor options; zero produces default behavior

Return Value

SPI_prepare_cursor has the same return conventions as SPI_prepare.

Notes

Useful bits to set in *cursorOptions* include `CURSOR_OPT_SCROLL`, `CURSOR_OPT_NO_SCROLL`, `CURSOR_OPT_FAST_PLAN`, `CURSOR_OPT_GENERIC_PLAN`, and `CURSOR_OPT_CUSTOM_PLAN`. Note in particular that `CURSOR_OPT_HOLD` is ignored.

SPI_prepare_extended

SPI_prepare_extended — prepare a statement, without executing it yet

Synopsis

```
SPIPlanPtr SPI_prepare_extended(const char * command,  
                               const SPIPrepareOptions * options)
```

Description

SPI_prepare_extended creates and returns a prepared statement for the specified command, but doesn't execute the command. This function is equivalent to SPI_prepare, with the addition that the caller can specify options to control the parsing of external parameter references, as well as other facets of query parsing and planning.

Arguments

const char * *command*

command string

const SPIPrepareOptions * *options*

struct containing optional arguments

Callers should always zero out the entire *options* struct, then fill whichever fields they want to set. This ensures forward compatibility of code, since any fields that are added to the struct in future will be defined to behave backwards-compatibly if they are zero. The currently available *options* fields are:

ParserSetupHook *parserSetup*

Parser hook setup function

void * *parserSetupArg*

pass-through argument for *parserSetup*

RawParseMode *parseMode*

mode for raw parsing; RAW_PARSE_DEFAULT (zero) produces default behavior

int *cursorOptions*

integer bit mask of cursor options; zero produces default behavior

Return Value

SPI_prepare_extended has the same return conventions as SPI_prepare.

SPI_prepare_params

SPI_prepare_params — prepare a statement, without executing it yet

Synopsis

```
SPIPlanPtr SPI_prepare_params(const char * command,
                             ParserSetupHook parserSetup,
                             void * parserSetupArg,
                             int cursorOptions)
```

Description

SPI_prepare_params creates and returns a prepared statement for the specified command, but doesn't execute the command. This function is equivalent to SPI_prepare_cursor, with the addition that the caller can specify parser hook functions to control the parsing of external parameter references.

This function is now deprecated in favor of SPI_prepare_extended.

Arguments

const char * *command*

command string

ParserSetupHook *parserSetup*

Parser hook setup function

void * *parserSetupArg*

pass-through argument for *parserSetup*

int *cursorOptions*

integer bit mask of cursor options; zero produces default behavior

Return Value

SPI_prepare_params has the same return conventions as SPI_prepare.

SPI_getargcount

SPI_getargcount — return the number of arguments needed by a statement prepared by SPI_prepare

Synopsis

```
int SPI_getargcount(SPIPlanPtr plan)
```

Description

SPI_getargcount returns the number of arguments needed to execute a statement prepared by SPI_prepare.

Arguments

SPIPlanPtr *plan*

prepared statement (returned by SPI_prepare)

Return Value

The count of expected arguments for the *plan*. If the *plan* is NULL or invalid, SPI_result is set to SPI_ERROR_ARGUMENT and -1 is returned.

SPI_getargtypeid

SPI_getargtypeid — return the data type OID for an argument of a statement prepared by SPI_prepare

Synopsis

```
Oid SPI_getargtypeid(SPIPlanPtr plan, int argIndex)
```

Description

SPI_getargtypeid returns the OID representing the type for the *argIndex*'th argument of a statement prepared by SPI_prepare. First argument is at index zero.

Arguments

SPIPlanPtr *plan*

prepared statement (returned by SPI_prepare)

int *argIndex*

zero based index of the argument

Return Value

The type OID of the argument at the given index. If the *plan* is NULL or invalid, or *argIndex* is less than 0 or not less than the number of arguments declared for the *plan*, SPI_result is set to SPI_ERROR_ARGUMENT and InvalidOid is returned.

SPI_is_cursor_plan

`SPI_is_cursor_plan` — return `true` if a statement prepared by `SPI_prepare` can be used with `SPI_cursor_open`

Synopsis

```
bool SPI_is_cursor_plan(SPIPlanPtr plan)
```

Description

`SPI_is_cursor_plan` returns `true` if a statement prepared by `SPI_prepare` can be passed as an argument to `SPI_cursor_open`, or `false` if that is not the case. The criteria are that the *plan* represents one single command and that this command returns tuples to the caller; for example, `SELECT` is allowed unless it contains an `INTO` clause, and `UPDATE` is allowed only if it contains a `RETURNING` clause.

Arguments

`SPIPlanPtr plan`

prepared statement (returned by `SPI_prepare`)

Return Value

`true` or `false` to indicate if the *plan* can produce a cursor or not, with `SPI_result` set to zero. If it is not possible to determine the answer (for example, if the *plan* is `NULL` or invalid, or if called when not connected to SPI), then `SPI_result` is set to a suitable error code and `false` is returned.

SPI_execute_plan

`SPI_execute_plan` — execute a statement prepared by `SPI_prepare`

Synopsis

```
int SPI_execute_plan(SPIPlanPtr plan, Datum * values, const char * nulls,
                    bool read_only, long count)
```

Description

`SPI_execute_plan` executes a statement prepared by `SPI_prepare` or one of its siblings. `read_only` and `count` have the same interpretation as in `SPI_execute`.

Arguments

`SPIPlanPtr plan`

prepared statement (returned by `SPI_prepare`)

`Datum * values`

An array of actual parameter values. Must have same length as the statement's number of arguments.

`const char * nulls`

An array describing which parameters are null. Must have same length as the statement's number of arguments.

If `nulls` is `NULL` then `SPI_execute_plan` assumes that no parameters are null. Otherwise, each entry of the `nulls` array should be ' ' if the corresponding parameter value is non-null, or 'n' if the corresponding parameter value is null. (In the latter case, the actual value in the corresponding `values` entry doesn't matter.) Note that `nulls` is not a text string, just an array: it does not need a '\0' terminator.

`bool read_only`

true for read-only execution

`long count`

maximum number of rows to return, or 0 for no limit

Return Value

The return value is the same as for `SPI_execute`, with the following additional possible error (negative) results:

`SPI_ERROR_ARGUMENT`

if `plan` is `NULL` or invalid, or `count` is less than 0

`SPI_ERROR_PARAM`

if `values` is `NULL` and `plan` was prepared with some parameters

`SPI_processed` and `SPI_tuptable` are set as in `SPI_execute` if successful.

SPI_execute_plan_extended

`SPI_execute_plan_extended` — execute a statement prepared by `SPI_prepare`

Synopsis

```
int SPI_execute_plan_extended(SPIPlanPtr plan,
                             const SPIExecuteOptions * options)
```

Description

`SPI_execute_plan_extended` executes a statement prepared by `SPI_prepare` or one of its siblings. This function is equivalent to `SPI_execute_plan`, except that information about the parameter values to be passed to the query is presented differently, and additional execution-controlling options can be passed.

Query parameter values are represented by a `ParamListInfo` struct, which is convenient for passing down values that are already available in that format. Dynamic parameter sets can also be used, via hook functions specified in `ParamListInfo`.

Also, instead of always accumulating the result tuples into a `SPI_tuptable` structure, tuples can be passed to a caller-supplied `DestReceiver` object as they are generated by the executor. This is particularly helpful for queries that might generate many tuples, since the data can be processed on-the-fly instead of being accumulated in memory.

Arguments

`SPIPlanPtr plan`

prepared statement (returned by `SPI_prepare`)

`const SPIExecuteOptions * options`

struct containing optional arguments

Callers should always zero out the entire `options` struct, then fill whichever fields they want to set. This ensures forward compatibility of code, since any fields that are added to the struct in future will be defined to behave backwards-compatibly if they are zero. The currently available `options` fields are:

`ParamListInfo params`

data structure containing query parameter types and values; NULL if none

`bool read_only`

true for read-only execution

`bool allow_nonatomic`

true allows non-atomic execution of CALL and DO statements (but this field is ignored unless the `SPI_OPT_NONATOMIC` flag was passed to `SPI_connect_ext`)

`bool must_return_tuples`

if true, raise error if the query is not of a kind that returns tuples (this does not forbid the case where it happens to return zero tuples)

`uint64 tcount`

maximum number of rows to return, or 0 for no limit

`DestReceiver * dest`

`DestReceiver` object that will receive any tuples emitted by the query; if NULL, result tuples are accumulated into a `SPI_tuptable` structure, as in `SPI_execute_plan`

ResourceOwner owner

The resource owner that will hold a reference count on the plan while it is executed. If NULL, CurrentResourceOwner is used. Ignored for non-saved plans, as SPI does not acquire reference counts on those.

Return Value

The return value is the same as for SPI_execute_plan.

When *options->dest* is NULL, SPI_processed and SPI_tuptable are set as in SPI_execute_plan. When *options->dest* is not NULL, SPI_processed is set to zero and SPI_tuptable is set to NULL. If a tuple count is required, the caller's DestReceiver object must calculate it.

SPI_execute_plan_with_paramlist

`SPI_execute_plan_with_paramlist` — execute a statement prepared by `SPI_prepare`

Synopsis

```
int SPI_execute_plan_with_paramlist(SPIPlanPtr plan,
                                   ParamListInfo params,
                                   bool read_only,
                                   long count)
```

Description

`SPI_execute_plan_with_paramlist` executes a statement prepared by `SPI_prepare`. This function is equivalent to `SPI_execute_plan` except that information about the parameter values to be passed to the query is presented differently. The `ParamListInfo` representation can be convenient for passing down values that are already available in that format. It also supports use of dynamic parameter sets via hook functions specified in `ParamListInfo`.

This function is now deprecated in favor of `SPI_execute_plan_extended`.

Arguments

`SPIPlanPtr plan`

prepared statement (returned by `SPI_prepare`)

`ParamListInfo params`

data structure containing parameter types and values; NULL if none

`bool read_only`

true for read-only execution

`long count`

maximum number of rows to return, or 0 for no limit

Return Value

The return value is the same as for `SPI_execute_plan`.

`SPI_processed` and `SPI_tuptable` are set as in `SPI_execute_plan` if successful.

SPI_execp

SPI_execp — execute a statement in read/write mode

Synopsis

```
int SPI_execp(SPIPlanPtr plan, Datum * values, const char * nulls, long count)
```

Description

SPI_execp is the same as SPI_execute_plan, with the latter's *read_only* parameter always taken as false.

Arguments

SPIPlanPtr *plan*

prepared statement (returned by SPI_prepare)

Datum * *values*

An array of actual parameter values. Must have same length as the statement's number of arguments.

const char * *nulls*

An array describing which parameters are null. Must have same length as the statement's number of arguments.

If *nulls* is NULL then SPI_execp assumes that no parameters are null. Otherwise, each entry of the *nulls* array should be ' ' if the corresponding parameter value is non-null, or 'n' if the corresponding parameter value is null. (In the latter case, the actual value in the corresponding *values* entry doesn't matter.) Note that *nulls* is not a text string, just an array: it does not need a '\0' terminator.

long *count*

maximum number of rows to return, or 0 for no limit

Return Value

See SPI_execute_plan.

SPI_processed and SPI_tuptable are set as in SPI_execute if successful.

SPI_cursor_open

`SPI_cursor_open` — set up a cursor using a statement created with `SPI_prepare`

Synopsis

```
Portal SPI_cursor_open(const char * name, SPIPlanPtr plan,
                      Datum * values, const char * nulls,
                      bool read_only)
```

Description

`SPI_cursor_open` sets up a cursor (internally, a portal) that will execute a statement prepared by `SPI_prepare`. The parameters have the same meanings as the corresponding parameters to `SPI_execute_plan`.

Using a cursor instead of executing the statement directly has two benefits. First, the result rows can be retrieved a few at a time, avoiding memory overrun for queries that return many rows. Second, a portal can outlive the current C function (it can, in fact, live to the end of the current transaction). Returning the portal name to the C function's caller provides a way of returning a row set as result.

The passed-in parameter data will be copied into the cursor's portal, so it can be freed while the cursor still exists.

Arguments

`const char * name`

name for portal, or `NULL` to let the system select a name

`SPIPlanPtr plan`

prepared statement (returned by `SPI_prepare`)

`Datum * values`

An array of actual parameter values. Must have same length as the statement's number of arguments.

`const char * nulls`

An array describing which parameters are null. Must have same length as the statement's number of arguments.

If `nulls` is `NULL` then `SPI_cursor_open` assumes that no parameters are null. Otherwise, each entry of the `nulls` array should be ' ' if the corresponding parameter value is non-null, or 'n' if the corresponding parameter value is null. (In the latter case, the actual value in the corresponding `values` entry doesn't matter.) Note that `nulls` is not a text string, just an array: it does not need a '\0' terminator.

`bool read_only`

true for read-only execution

Return Value

Pointer to portal containing the cursor. Note there is no error return convention; any error will be reported via `elog`.

SPI_cursor_open_with_args

`SPI_cursor_open_with_args` — set up a cursor using a query and parameters

Synopsis

```
Portal SPI_cursor_open_with_args(const char *name,
                                const char *command,
                                int nargs, Oid *argtypes,
                                Datum *values, const char *nulls,
                                bool read_only, int cursorOptions)
```

Description

`SPI_cursor_open_with_args` sets up a cursor (internally, a portal) that will execute the specified query. Most of the parameters have the same meanings as the corresponding parameters to `SPI_prepare_cursor` and `SPI_cursor_open`.

For one-time query execution, this function should be preferred over `SPI_prepare_cursor` followed by `SPI_cursor_open`. If the same command is to be executed with many different parameters, either method might be faster, depending on the cost of re-planning versus the benefit of custom plans.

The passed-in parameter data will be copied into the cursor's portal, so it can be freed while the cursor still exists.

This function is now deprecated in favor of `SPI_cursor_parse_open`, which provides equivalent functionality using a more modern API for handling query parameters.

Arguments

`const char * name`

name for portal, or `NULL` to let the system select a name

`const char * command`

command string

`int nargs`

number of input parameters (\$1, \$2, etc.)

`Oid * argtypes`

an array of length `nargs`, containing the OIDs of the data types of the parameters

`Datum * values`

an array of length `nargs`, containing the actual parameter values

`const char * nulls`

an array of length `nargs`, describing which parameters are null

If `nulls` is `NULL` then `SPI_cursor_open_with_args` assumes that no parameters are null. Otherwise, each entry of the `nulls` array should be ' ' if the corresponding parameter value is non-null, or 'n' if the corresponding parameter value is null. (In the latter case, the actual value in the corresponding `values` entry doesn't matter.) Note that `nulls` is not a text string, just an array: it does not need a '\0' terminator.

`bool read_only`

true for read-only execution

`int cursorOptions`

integer bit mask of cursor options; zero produces default behavior

Return Value

Pointer to portal containing the cursor. Note there is no error return convention; any error will be reported via `elog`.

SPI_cursor_open_with_paramlist

SPI_cursor_open_with_paramlist — set up a cursor using parameters

Synopsis

```
Portal SPI_cursor_open_with_paramlist(const char *name,  
                                     SPIPlanPtr plan,  
                                     ParamListInfo params,  
                                     bool read_only)
```

Description

SPI_cursor_open_with_paramlist sets up a cursor (internally, a portal) that will execute a statement prepared by SPI_prepare. This function is equivalent to SPI_cursor_open except that information about the parameter values to be passed to the query is presented differently. The ParamListInfo representation can be convenient for passing down values that are already available in that format. It also supports use of dynamic parameter sets via hook functions specified in ParamListInfo.

The passed-in parameter data will be copied into the cursor's portal, so it can be freed while the cursor still exists.

Arguments

const char * name

name for portal, or NULL to let the system select a name

SPIPlanPtr plan

prepared statement (returned by SPI_prepare)

ParamListInfo params

data structure containing parameter types and values; NULL if none

bool read_only

true for read-only execution

Return Value

Pointer to portal containing the cursor. Note there is no error return convention; any error will be reported via elog.

SPI_cursor_parse_open

SPI_cursor_parse_open — set up a cursor using a query string and parameters

Synopsis

```
Portal SPI_cursor_parse_open(const char *name,
                             const char *command,
                             const SPIParseOpenOptions * options)
```

Description

SPI_cursor_parse_open sets up a cursor (internally, a portal) that will execute the specified query string. This is comparable to SPI_prepare_cursor followed by SPI_cursor_open_with_paramlist, except that parameter references within the query string are handled entirely by supplying a ParamListInfo object.

For one-time query execution, this function should be preferred over SPI_prepare_cursor followed by SPI_cursor_open_with_paramlist. If the same command is to be executed with many different parameters, either method might be faster, depending on the cost of re-planning versus the benefit of custom plans.

The *options->params* object should normally mark each parameter with the PARAM_FLAG_CONST flag, since a one-shot plan is always used for the query.

The passed-in parameter data will be copied into the cursor's portal, so it can be freed while the cursor still exists.

Arguments

const char * *name*
name for portal, or NULL to let the system select a name

const char * *command*
command string

const SPIParseOpenOptions * *options*
struct containing optional arguments

Callers should always zero out the entire *options* struct, then fill whichever fields they want to set. This ensures forward compatibility of code, since any fields that are added to the struct in future will be defined to behave backwards-compatibly if they are zero. The currently available *options* fields are:

ParamListInfo *params*
data structure containing query parameter types and values; NULL if none

int *cursorOptions*
integer bit mask of cursor options; zero produces default behavior

bool *read_only*
true for read-only execution

Return Value

Pointer to portal containing the cursor. Note there is no error return convention; any error will be reported via elog.

SPI_cursor_find

SPI_cursor_find — find an existing cursor by name

Synopsis

```
Portal SPI_cursor_find(const char * name)
```

Description

SPI_cursor_find finds an existing portal by name. This is primarily useful to resolve a cursor name returned as text by some other function.

Arguments

```
const char * name
```

name of the portal

Return Value

pointer to the portal with the specified name, or `NULL` if none was found

Notes

Beware that this function can return a `Portal` object that does not have cursor-like properties; for example it might not return tuples. If you simply pass the `Portal` pointer to other SPI functions, they can defend themselves against such cases, but caution is appropriate when directly inspecting the `Portal`.

SPI_cursor_fetch

SPI_cursor_fetch — fetch some rows from a cursor

Synopsis

```
void SPI_cursor_fetch(Portal portal, bool forward, long count)
```

Description

SPI_cursor_fetch fetches some rows from a cursor. This is equivalent to a subset of the SQL command FETCH (see SPI_scroll_cursor_fetch for more functionality).

Arguments

Portal *portal*

portal containing the cursor

bool *forward*

true for fetch forward, false for fetch backward

long *count*

maximum number of rows to fetch

Return Value

SPI_processed and SPI_tuptable are set as in SPI_execute if successful.

Notes

Fetching backward may fail if the cursor's plan was not created with the CURSOR_OPT_SCROLL option.

SPI_cursor_move

SPI_cursor_move — move a cursor

Synopsis

```
void SPI_cursor_move(Portal portal, bool forward, long count)
```

Description

SPI_cursor_move skips over some number of rows in a cursor. This is equivalent to a subset of the SQL command `MOVE` (see SPI_scroll_cursor_move for more functionality).

Arguments

Portal *portal*

portal containing the cursor

bool *forward*

true for move forward, false for move backward

long *count*

maximum number of rows to move

Notes

Moving backward may fail if the cursor's plan was not created with the `CURSOR_OPT_SCROLL` option.

SPI_scroll_cursor_fetch

SPI_scroll_cursor_fetch — fetch some rows from a cursor

Synopsis

```
void SPI_scroll_cursor_fetch(Portal portal, FetchDirection direction,
                             long count)
```

Description

SPI_scroll_cursor_fetch fetches some rows from a cursor. This is equivalent to the SQL command `FETCH`.

Arguments

Portal *portal*

portal containing the cursor

FetchDirection *direction*

one of `FETCH_FORWARD`, `FETCH_BACKWARD`, `FETCH_ABSOLUTE` or `FETCH_RELATIVE`

long *count*

number of rows to fetch for `FETCH_FORWARD` or `FETCH_BACKWARD`; absolute row number to fetch for `FETCH_ABSOLUTE`; or relative row number to fetch for `FETCH_RELATIVE`

Return Value

SPI_processed and SPI_tuptable are set as in SPI_execute if successful.

Notes

See the SQL [FETCH](#) command for details of the interpretation of the *direction* and *count* parameters.

Direction values other than `FETCH_FORWARD` may fail if the cursor's plan was not created with the `CURSOR_OPT_SCROLL` option.

SPI_scroll_cursor_move

SPI_scroll_cursor_move — move a cursor

Synopsis

```
void SPI_scroll_cursor_move(Portal portal, FetchDirection direction,
                           long count)
```

Description

SPI_scroll_cursor_move skips over some number of rows in a cursor. This is equivalent to the SQL command MOVE.

Arguments

Portal *portal*

portal containing the cursor

FetchDirection *direction*

one of FETCH_FORWARD, FETCH_BACKWARD, FETCH_ABSOLUTE or FETCH_RELATIVE

long *count*

number of rows to move for FETCH_FORWARD or FETCH_BACKWARD; absolute row number to move to for FETCH_ABSOLUTE; or relative row number to move to for FETCH_RELATIVE

Return Value

SPI_processed is set as in SPI_execute if successful. SPI_tuptable is set to NULL, since no rows are returned by this function.

Notes

See the SQL [FETCH](#) command for details of the interpretation of the *direction* and *count* parameters.

Direction values other than FETCH_FORWARD may fail if the cursor's plan was not created with the CURSOR_OPT_SCROLL option.

SPI_cursor_close

SPI_cursor_close — close a cursor

Synopsis

```
void SPI_cursor_close(Portal portal)
```

Description

SPI_cursor_close closes a previously created cursor and releases its portal storage.

All open cursors are closed automatically at the end of a transaction. SPI_cursor_close need only be invoked if it is desirable to release resources sooner.

Arguments

Portal *portal*

portal containing the cursor

SPI_keepplan

SPI_keepplan — save a prepared statement

Synopsis

```
int SPI_keepplan(SPIPlanPtr plan)
```

Description

SPI_keepplan saves a passed statement (prepared by SPI_prepare) so that it will not be freed by SPI_finish nor by the transaction manager. This gives you the ability to reuse prepared statements in the subsequent invocations of your C function in the current session.

Arguments

SPIPlanPtr *plan*

the prepared statement to be saved

Return Value

0 on success; SPI_ERROR_ARGUMENT if *plan* is NULL or invalid

Notes

The passed-in statement is relocated to permanent storage by means of pointer adjustment (no data copying is required). If you later wish to delete it, use SPI_freeplan on it.

SPI_saveplan

SPI_saveplan — save a prepared statement

Synopsis

```
SPIPlanPtr SPI_saveplan(SPIPlanPtr plan)
```

Description

SPI_saveplan copies a passed statement (prepared by SPI_prepare) into memory that will not be freed by SPI_finish nor by the transaction manager, and returns a pointer to the copied statement. This gives you the ability to reuse prepared statements in the subsequent invocations of your C function in the current session.

Arguments

SPIPlanPtr *plan*

the prepared statement to be saved

Return Value

Pointer to the copied statement; or NULL if unsuccessful. On error, SPI_result is set thus:

SPI_ERROR_ARGUMENT

if *plan* is NULL or invalid

SPI_ERROR_UNCONNECTED

if called from an unconnected C function

Notes

The originally passed-in statement is not freed, so you might wish to do SPI_freeplan on it to avoid leaking memory until SPI_finish.

In most cases, SPI_keepplan is preferred to this function, since it accomplishes largely the same result without needing to physically copy the prepared statement's data structures.

SPI_register_relation

SPI_register_relation — make an ephemeral named relation available by name in SPI queries

Synopsis

```
int SPI_register_relation(EphemeralNamedRelation enr)
```

Description

SPI_register_relation makes an ephemeral named relation, with associated information, available to queries planned and executed through the current SPI connection.

Arguments

EphemeralNamedRelation *enr*

the ephemeral named relation registry entry

Return Value

If the execution of the command was successful then the following (nonnegative) value will be returned:

SPI_OK_REL_REGISTER

if the relation has been successfully registered by name

On error, one of the following negative values is returned:

SPI_ERROR_ARGUMENT

if *enr* is NULL or its `name` field is NULL

SPI_ERROR_UNCONNECTED

if called from an unconnected C function

SPI_ERROR_REL_DUPLICATE

if the name specified in the `name` field of *enr* is already registered for this connection

SPI_unregister_relation

SPI_unregister_relation — remove an ephemeral named relation from the registry

Synopsis

```
int SPI_unregister_relation(const char * name)
```

Description

SPI_unregister_relation removes an ephemeral named relation from the registry for the current connection.

Arguments

const char * name

the relation registry entry name

Return Value

If the execution of the command was successful then the following (nonnegative) value will be returned:

SPI_OK_REL_UNREGISTER

if the tuplestore has been successfully removed from the registry

On error, one of the following negative values is returned:

SPI_ERROR_ARGUMENT

if *name* is NULL

SPI_ERROR_UNCONNECTED

if called from an unconnected C function

SPI_ERROR_REL_NOT_FOUND

if *name* is not found in the registry for the current connection

SPI_register_trigger_data

SPI_register_trigger_data — make ephemeral trigger data available in SPI queries

Synopsis

```
int SPI_register_trigger_data(TriggerData *tdata)
```

Description

SPI_register_trigger_data makes any ephemeral relations captured by a trigger available to queries planned and executed through the current SPI connection. Currently, this means the transition tables captured by an AFTER trigger defined with a REFERENCING OLD/NEW TABLE AS ... clause. This function should be called by a PL trigger handler function after connecting.

Arguments

TriggerData *tdata

the TriggerData object passed to a trigger handler function as fcinfo->context

Return Value

If the execution of the command was successful then the following (nonnegative) value will be returned:

SPI_OK_TD_REGISTER

if the captured trigger data (if any) has been successfully registered

On error, one of the following negative values is returned:

SPI_ERROR_ARGUMENT

if tdata is NULL

SPI_ERROR_UNCONNECTED

if called from an unconnected C function

SPI_ERROR_REL_DUPLICATE

if the name of any trigger data transient relation is already registered for this connection

48.2. Interface Support Functions

The functions described here provide an interface for extracting information from result sets returned by SPI_execute and other SPI functions.

All functions described in this section can be used by both connected and unconnected C functions.

SPI_fname

SPI_fname — determine the column name for the specified column number

Synopsis

```
char * SPI_fname(TupleDesc rowdesc, int colnumber)
```

Description

SPI_fname returns a copy of the column name of the specified column. (You can use `pfree` to release the copy of the name when you don't need it anymore.)

Arguments

`TupleDesc rowdesc`

input row description

`int colnumber`

column number (count starts at 1)

Return Value

The column name; NULL if `colnumber` is out of range. SPI_result set to SPI_ERROR_NOATTRIBUTE on error.

SPI_fnumber

SPI_fnumber — determine the column number for the specified column name

Synopsis

```
int SPI_fnumber(TupleDesc rowdesc, const char * colname)
```

Description

SPI_fnumber returns the column number for the column with the specified name.

If *colname* refers to a system column (e.g., *ctid*) then the appropriate negative column number will be returned. The caller should be careful to test the return value for exact equality to SPI_ERROR_NOATTRIBUTE to detect an error; testing the result for less than or equal to 0 is not correct unless system columns should be rejected.

Arguments

```
TupleDesc rowdesc  
    input row description  
  
const char * colname  
    column name
```

Return Value

Column number (count starts at 1 for user-defined columns), or SPI_ERROR_NOATTRIBUTE if the named column was not found.

SPI_getvalue

SPI_getvalue — return the string value of the specified column

Synopsis

```
char * SPI_getvalue(HeapTuple row, TupleDesc rowdesc, int colnumber)
```

Description

SPI_getvalue returns the string representation of the value of the specified column.

The result is returned in memory allocated using `palloc`. (You can use `pfree` to release the memory when you don't need it anymore.)

Arguments

HeapTuple row

input row to be examined

TupleDesc rowdesc

input row description

int colnumber

column number (count starts at 1)

Return Value

Column value, or `NULL` if the column is null, *colnumber* is out of range (SPI_result is set to SPI_ERROR_NOATTRIBUTE), or no output function is available (SPI_result is set to SPI_ERROR_NOOUTFUNC).

SPI_getbinval

SPI_getbinval — return the binary value of the specified column

Synopsis

```
Datum SPI_getbinval(HeapTuple row, TupleDesc rowdesc, int colnumber,  
                    bool * isnull)
```

Description

SPI_getbinval returns the value of the specified column in the internal form (as type Datum).

This function does not allocate new space for the datum. In the case of a pass-by-reference data type, the return value will be a pointer into the passed row.

Arguments

HeapTuple row

input row to be examined

TupleDesc rowdesc

input row description

int colnumber

column number (count starts at 1)

bool * isnull

flag for a null value in the column

Return Value

The binary value of the column is returned. The variable pointed to by *isnull* is set to true if the column is null, else to false.

SPI_result is set to SPI_ERROR_NOATTRIBUTE on error.

SPI_gettype

SPI_gettype — return the data type name of the specified column

Synopsis

```
char * SPI_gettype(TupleDesc rowdesc, int colnumber)
```

Description

SPI_gettype returns a copy of the data type name of the specified column. (You can use `pfree` to release the copy of the name when you don't need it anymore.)

Arguments

`TupleDesc rowdesc`

input row description

`int colnumber`

column number (count starts at 1)

Return Value

The data type name of the specified column, or `NULL` on error. `SPI_result` is set to `SPI_ERROR_NOATTRIBUTE` on error.

SPI_gettypeid

SPI_gettypeid — return the data type OID of the specified column

Synopsis

```
Oid SPI_gettypeid(TupleDesc rowdesc, int colnumber)
```

Description

SPI_gettypeid returns the OID of the data type of the specified column.

Arguments

TupleDesc rowdesc

input row description

int colnumber

column number (count starts at 1)

Return Value

The OID of the data type of the specified column or `InvalidOid` on error. On error, `SPI_result` is set to `SPI_ERROR_NOATTRIBUTE`.

SPI_getrelname

SPI_getrelname — return the name of the specified relation

Synopsis

```
char * SPI_getrelname(Relation rel)
```

Description

SPI_getrelname returns a copy of the name of the specified relation. (You can use pfree to release the copy of the name when you don't need it anymore.)

Arguments

Relation *rel*
input relation

Return Value

The name of the specified relation.

SPI_getnspname

SPI_getnspname — return the namespace of the specified relation

Synopsis

```
char * SPI_getnspname(Relation rel)
```

Description

SPI_getnspname returns a copy of the name of the namespace that the specified `Relation` belongs to. This is equivalent to the relation's schema. You should `pfree` the return value of this function when you are finished with it.

Arguments

`Relation rel`
input relation

Return Value

The name of the specified relation's namespace.

SPI_result_code_string

SPI_result_code_string — return error code as string

Synopsis

```
const char * SPI_result_code_string(int code);
```

Description

SPI_result_code_string returns a string representation of the result code returned by various SPI functions or stored in SPI_result.

Arguments

int code
result code

Return Value

A string representation of the result code.

48.3. Memory Management

Postgres Pro allocates memory within *memory contexts*, which provide a convenient method of managing allocations made in many different places that need to live for differing amounts of time. Destroying a context releases all the memory that was allocated in it. Thus, it is not necessary to keep track of individual objects to avoid memory leaks; instead only a relatively small number of contexts have to be managed. `palloc` and related functions allocate memory from the “current” context.

SPI_connect creates a new memory context and makes it current. SPI_finish restores the previous current memory context and destroys the context created by SPI_connect. These actions ensure that transient memory allocations made inside your C function are reclaimed at C function exit, avoiding memory leakage.

However, if your C function needs to return an object in allocated memory (such as a value of a pass-by-reference data type), you cannot allocate that memory using `palloc`, at least not while you are connected to SPI. If you try, the object will be deallocated by SPI_finish, and your C function will not work reliably. To solve this problem, use `SPI_palloc` to allocate memory for your return object. `SPI_palloc` allocates memory in the “upper executor context”, that is, the memory context that was current when SPI_connect was called, which is precisely the right context for a value returned from your C function. Several of the other utility functions described in this section also return objects created in the upper executor context.

When SPI_connect is called, the private context of the C function, which is created by SPI_connect, is made the current context. All allocations made by `palloc`, `repalloc`, or SPI utility functions (except as described in this section) are made in this context. When a C function disconnects from the SPI manager (via SPI_finish) the current context is restored to the upper executor context, and all allocations made in the C function memory context are freed and cannot be used any more.

SPI_palloc

SPI_palloc — allocate memory in the upper executor context

Synopsis

```
void * SPI_palloc(Size size)
```

Description

SPI_palloc allocates memory in the upper executor context.

This function can only be used while connected to SPI. Otherwise, it throws an error.

Arguments

Size *size*

size in bytes of storage to allocate

Return Value

pointer to new storage space of the specified size

SPI_realloc

SPI_realloc — reallocate memory in the upper executor context

Synopsis

```
void * SPI_realloc(void * pointer, Size size)
```

Description

SPI_realloc changes the size of a memory segment previously allocated using SPI_palloc.

This function is no longer different from plain realloc. It's kept just for backward compatibility of existing code.

Arguments

`void * pointer`

pointer to existing storage to change

`Size size`

size in bytes of storage to allocate

Return Value

pointer to new storage space of specified size with the contents copied from the existing area

SPI_pfree

SPI_pfree — free memory in the upper executor context

Synopsis

```
void SPI_pfree(void * pointer)
```

Description

SPI_pfree frees memory previously allocated using SPI_palloc or SPI_realloc.

This function is no longer different from plain pfree. It's kept just for backward compatibility of existing code.

Arguments

```
void * pointer
```

pointer to existing storage to free

SPI_copytuple

SPI_copytuple — make a copy of a row in the upper executor context

Synopsis

```
HeapTuple SPI_copytuple(HeapTuple row)
```

Description

SPI_copytuple makes a copy of a row in the upper executor context. This is normally used to return a modified row from a trigger. In a function declared to return a composite type, use SPI_returntuple instead.

This function can only be used while connected to SPI. Otherwise, it returns NULL and sets SPI_result to SPI_ERROR_UNCONNECTED.

Arguments

HeapTuple row
row to be copied

Return Value

the copied row, or NULL on error (see SPI_result for an error indication)

SPI_returntuple

SPI_returntuple — prepare to return a tuple as a Datum

Synopsis

```
HeapTupleHeader SPI_returntuple(HeapTuple row, TupleDesc rowdesc)
```

Description

SPI_returntuple makes a copy of a row in the upper executor context, returning it in the form of a row type Datum. The returned pointer need only be converted to Datum via PointerGetDatum before returning.

This function can only be used while connected to SPI. Otherwise, it returns NULL and sets SPI_result to SPI_ERROR_UNCONNECTED.

Note that this should be used for functions that are declared to return composite types. It is not used for triggers; use SPI_copytuple for returning a modified row in a trigger.

Arguments

HeapTuple row

row to be copied

TupleDesc rowdesc

descriptor for row (pass the same descriptor each time for most effective caching)

Return Value

HeapTupleHeader pointing to copied row, or NULL on error (see SPI_result for an error indication)

SPI_modifytuple

SPI_modifytuple — create a row by replacing selected fields of a given row

Synopsis

```
HeapTuple SPI_modifytuple(Relation rel, HeapTuple row, int ncols,
                          int * colnum, Datum * values, const char * nulls)
```

Description

SPI_modifytuple creates a new row by substituting new values for selected columns, copying the original row's columns at other positions. The input row is not modified. The new row is returned in the upper executor context.

This function can only be used while connected to SPI. Otherwise, it returns NULL and sets SPI_result to SPI_ERROR_UNCONNECTED.

Arguments

Relation *rel*

Used only as the source of the row descriptor for the row. (Passing a relation rather than a row descriptor is a misfeature.)

HeapTuple *row*

row to be modified

int *ncols*

number of columns to be changed

int * *colnum*

an array of length *ncols*, containing the numbers of the columns that are to be changed (column numbers start at 1)

Datum * *values*

an array of length *ncols*, containing the new values for the specified columns

const char * *nulls*

an array of length *ncols*, describing which new values are null

If *nulls* is NULL then SPI_modifytuple assumes that no new values are null. Otherwise, each entry of the *nulls* array should be ' ' if the corresponding new value is non-null, or 'n' if the corresponding new value is null. (In the latter case, the actual value in the corresponding *values* entry doesn't matter.) Note that *nulls* is not a text string, just an array: it does not need a '\0' terminator.

Return Value

new row with modifications, allocated in the upper executor context, or NULL on error (see SPI_result for an error indication)

On error, SPI_result is set as follows:

SPI_ERROR_ARGUMENT

if *rel* is NULL, or if *row* is NULL, or if *ncols* is less than or equal to 0, or if *colnum* is NULL, or if *values* is NULL.

`SPI_ERROR_NOATTRIBUTE`

if *colnum* contains an invalid column number (less than or equal to 0 or greater than the number of columns in *row*)

`SPI_ERROR_UNCONNECTED`

if SPI is not active

SPI_freetuple

SPI_freetuple — free a row allocated in the upper executor context

Synopsis

```
void SPI_freetuple(HeapTuple row)
```

Description

SPI_freetuple frees a row previously allocated in the upper executor context.

This function is no longer different from plain heap_freetuple. It's kept just for backward compatibility of existing code.

Arguments

HeapTuple row
row to free

SPI_freetuptable

SPI_freetuptable — free a row set created by SPI_execute or a similar function

Synopsis

```
void SPI_freetuptable(SPITupleTable * tuptable)
```

Description

SPI_freetuptable frees a row set created by a prior SPI command execution function, such as SPI_execute. Therefore, this function is often called with the global variable SPI_tuptable as argument.

This function is useful if an SPI-using C function needs to execute multiple commands and does not want to keep the results of earlier commands around until it ends. Note that any unfreed row sets will be freed anyway at SPI_finish. Also, if a subtransaction is started and then aborted within execution of an SPI-using C function, SPI automatically frees any row sets created while the subtransaction was running.

Beginning in PostgreSQL 9.3, SPI_freetuptable contains guard logic to protect against duplicate deletion requests for the same row set. In previous releases, duplicate deletions would lead to crashes.

Arguments

```
SPITupleTable * tuptable
```

pointer to row set to free, or NULL to do nothing

SPI_freeplan

SPI_freeplan — free a previously saved prepared statement

Synopsis

```
int SPI_freeplan(SPIPlanPtr plan)
```

Description

SPI_freeplan releases a prepared statement previously returned by SPI_prepare or saved by SPI_keepplan or SPI_saveplan.

Arguments

SPIPlanPtr *plan*

pointer to statement to free

Return Value

0 on success; SPI_ERROR_ARGUMENT if *plan* is NULL or invalid

48.4. Transaction Management

It is not possible to run transaction control commands such as COMMIT and ROLLBACK through SPI functions such as SPI_execute. There are, however, separate interface functions that allow transaction control through SPI.

It is not generally safe and sensible to start and end transactions in arbitrary user-defined SQL-callable functions without taking into account the context in which they are called. For example, a transaction boundary in the middle of a function that is part of a complex SQL expression that is part of some SQL command will probably result in obscure internal errors or crashes. The interface functions presented here are primarily intended to be used by procedural language implementations to support transaction management in SQL-level procedures that are invoked by the CALL command, taking the context of the CALL invocation into account. SPI-using procedures implemented in C can implement the same logic, but the details of that are beyond the scope of this documentation.

SPI_commit

`SPI_commit`, `SPI_commit_and_chain` — commit the current transaction

Synopsis

```
void SPI_commit(void)
void SPI_commit_and_chain(void)
```

Description

`SPI_commit` commits the current transaction. It is approximately equivalent to running the SQL command `COMMIT`. After the transaction is committed, a new transaction is automatically started using default transaction characteristics, so that the caller can continue using SPI facilities. If there is a failure during commit, the current transaction is instead rolled back and a new transaction is started, after which the error is thrown in the usual way.

`SPI_commit_and_chain` is the same, but the new transaction is started with the same transaction characteristics as the just finished one, like with the SQL command `COMMIT AND CHAIN`.

These functions can only be executed if the SPI connection has been set as nonatomic in the call to `SPI_connect_ext`.

SPI_rollback

SPI_rollback, SPI_rollback_and_chain — abort the current transaction

Synopsis

```
void SPI_rollback(void)
void SPI_rollback_and_chain(void)
```

Description

SPI_rollback rolls back the current transaction. It is approximately equivalent to running the SQL command `ROLLBACK`. After the transaction is rolled back, a new transaction is automatically started using default transaction characteristics, so that the caller can continue using SPI facilities.

SPI_rollback_and_chain is the same, but the new transaction is started with the same transaction characteristics as the just finished one, like with the SQL command `ROLLBACK AND CHAIN`.

These functions can only be executed if the SPI connection has been set as nonatomic in the call to SPI_connect_ext.

SPI_start_transaction

SPI_start_transaction — obsolete function

Synopsis

```
void SPI_start_transaction(void)
```

Description

SPI_start_transaction does nothing, and exists only for code compatibility with earlier PostgreSQL releases. It used to be required after calling SPI_commit or SPI_rollback, but now those functions start a new transaction automatically.

48.5. Visibility of Data Changes

The following rules govern the visibility of data changes in functions that use SPI (or any other C function):

- During the execution of an SQL command, any data changes made by the command are invisible to the command itself. For example, in:

```
INSERT INTO a SELECT * FROM a;
```

the inserted rows are invisible to the SELECT part.
- Changes made by a command C are visible to all commands that are started after C, no matter whether they are started inside C (during the execution of C) or after C is done.
- Commands executed via SPI inside a function called by an SQL command (either an ordinary function or a trigger) follow one or the other of the above rules depending on the read/write flag passed to SPI. Commands executed in read-only mode follow the first rule: they cannot see changes of the calling command. Commands executed in read-write mode follow the second rule: they can see all changes made so far.
- All standard procedural languages set the SPI read-write mode depending on the volatility attribute of the function. Commands of STABLE and IMMUTABLE functions are done in read-only mode, while commands of VOLATILE functions are done in read-write mode. While authors of C functions are able to violate this convention, it's unlikely to be a good idea to do so.

The next section contains an example that illustrates the application of these rules.

48.6. Examples

This section contains a very simple example of SPI usage. The C function `execq` takes an SQL command as its first argument and a row count as its second, executes the command using `SPI_exec` and returns the number of rows that were processed by the command.

```
#include "postgres.h"

#include "executor/spi.h"
#include "utils/builtins.h"

PG_MODULE_MAGIC;

PG_FUNCTION_INFO_V1(execq);

Datum
execq(PG_FUNCTION_ARGS)
{
    char *command;
    int cnt;
    int ret;
```



```
uint64 proc;

/* Convert given text object to a C string */
command = text_to_cstring(PG_GETARG_TEXT_PP(0));
cnt = PG_GETARG_INT32(1);

SPI_connect();

ret = SPI_exec(command, cnt);

proc = SPI_processed;

/*
 * If some rows were fetched, print them via elog(INFO).
 */
if (ret > 0 && SPI_tuptable != NULL)
{
    SPITupleTable *tuptable = SPI_tuptable;
    TupleDesc tupdesc = tuptable->tupdesc;
    char buf[8192];
    uint64 j;

    for (j = 0; j < tuptable->numvals; j++)
    {
        HeapTuple tuple = tuptable->vals[j];
        int i;

        for (i = 1, buf[0] = 0; i <= tupdesc->natts; i++)
            snprintf(buf + strlen(buf), sizeof(buf) - strlen(buf), " %s%s",
                     SPI_getvalue(tuple, tupdesc, i),
                     (i == tupdesc->natts) ? " " : " |");
        elog(INFO, "EXECQ: %s", buf);
    }
}

SPI_finish();
pfree(command);

PG_RETURN_INT64(proc);
}
```

This is how you declare the function after having compiled it into a shared library (details are in [Section 39.10.5](#)):

```
CREATE FUNCTION execq(text, integer) RETURNS int8
AS 'filename'
LANGUAGE C STRICT;
```

Here is a sample session:

```
=> SELECT execq('CREATE TABLE a (x integer)', 0);
 execq
-----
      0
(1 row)

=> INSERT INTO a VALUES (execq('INSERT INTO a VALUES (0)', 0));
INSERT 0 1
=> SELECT execq('SELECT * FROM a', 0);
INFO:  EXECQ:  0      -- inserted by execq
```

```

INFO: EXECQ: 1      -- returned by execq and inserted by upper INSERT

  execq
-----
      2
(1 row)

=> SELECT execq('INSERT INTO a SELECT x + 2 FROM a RETURNING *', 1);
INFO: EXECQ: 2      -- 0 + 2, then execution was stopped by count
  execq
-----
      1
(1 row)

=> SELECT execq('SELECT * FROM a', 10);
INFO: EXECQ: 0
INFO: EXECQ: 1
INFO: EXECQ: 2

  execq
-----
      3      -- 10 is the max value only, 3 is the real number of rows
(1 row)

=> SELECT execq('INSERT INTO a SELECT x + 10 FROM a', 1);
  execq
-----
      3      -- all rows processed; count does not stop it, because nothing is
returned
(1 row)

=> SELECT * FROM a;
  x
----
  0
  1
  2
10
11
12
(6 rows)

=> DELETE FROM a;
DELETE 6
=> INSERT INTO a VALUES (execq('SELECT * FROM a', 0) + 1);
INSERT 0 1
=> SELECT * FROM a;
  x
---
  1      -- 0 (no rows in a) + 1
(1 row)

=> INSERT INTO a VALUES (execq('SELECT * FROM a', 0) + 1);
INFO: EXECQ: 1
INSERT 0 1
=> SELECT * FROM a;
  x
---
```

```

1
2          -- 1 (there was one row in a) + 1
(2 rows)

-- This demonstrates the data changes visibility rule.
-- execq is called twice and sees different numbers of rows each time:

=> INSERT INTO a SELECT execq('SELECT * FROM a', 0) * x FROM a;
INFO: EXECQ: 1      -- results from first execq
INFO: EXECQ: 2
INFO: EXECQ: 1      -- results from second execq
INFO: EXECQ: 2
INFO: EXECQ: 2
INSERT 0 2
=> SELECT * FROM a;
 x
---
1
2
2          -- 2 rows * 1 (x in first row)
6          -- 3 rows (2 + 1 just inserted) * 2 (x in second row)
(4 rows)

```

Chapter 49. Background Worker Processes

Postgres Pro can be extended to run user-supplied code in separate processes. Such processes are started, stopped and monitored by `postgres`, which permits them to have a lifetime closely linked to the server's status. These processes are attached to Postgres Pro's shared memory area and have the option to connect to databases internally; they can also run multiple transactions serially, just like a regular client-connected server process. Also, by linking to `libpq` they can connect to the server and behave like a regular client application.

Warning

There are considerable robustness and security risks in using background worker processes because, being written in the `C` language, they have unrestricted access to data. Administrators wishing to enable modules that include background worker processes should exercise extreme caution. Only carefully audited modules should be permitted to run background worker processes.

Background workers can be initialized at the time that Postgres Pro is started by including the module name in `shared_preload_libraries`. A module wishing to run a background worker can register it by calling `RegisterBackgroundWorker (BackgroundWorker *worker)` from its `_PG_init()` function. Background workers can also be started after the system is up and running by calling `RegisterDynamicBackgroundWorker (BackgroundWorker *worker, BackgroundWorkerHandle **handle)`. Unlike `RegisterBackgroundWorker`, which can only be called from within the postmaster process, `RegisterDynamicBackgroundWorker` must be called from a regular backend or another background worker.

The structure `BackgroundWorker` is defined thus:

```
typedef void (*bgworker_main_type) (Datum main_arg);
typedef struct BackgroundWorker
{
    char          bgw_name[BGW_MAXLEN];
    char          bgw_type[BGW_MAXLEN];
    int           bgw_flags;
    BgWorkerStartTime bgw_start_time;
    int           bgw_restart_time; /* in seconds, or BGW_NEVER_RESTART */
    char          bgw_library_name[MAXPGPATH];
    char          bgw_function_name[BGW_MAXLEN];
    Datum         bgw_main_arg;
    char          bgw_extra[BGW_EXTRALEN];
    pid_t         bgw_notify_pid;
} BackgroundWorker;
```

`bgw_name` and `bgw_type` are strings to be used in log messages, process listings and similar contexts. `bgw_type` should be the same for all background workers of the same type, so that it is possible to group such workers in a process listing, for example. `bgw_name` on the other hand can contain additional information about the specific process. (Typically, the string for `bgw_name` will contain the type somehow, but that is not strictly required.)

`bgw_flags` is a bitwise-or'd bit mask indicating the capabilities that the module wants. Possible values are:

`BGWORKER_SHMEM_ACCESS`

Requests shared memory access. This flag is required.

`BGWORKER_BACKEND_DATABASE_CONNECTION`

Requests the ability to establish a database connection through which it can later run transactions and queries. A background worker using `BGWORKER_BACKEND_DATABASE_CONNECTION` to connect to a database must also attach shared memory using `BGWORKER_SHMEM_ACCESS`, or worker start-up will fail.

`bgw_start_time` is the server state during which `postgres` should start the process; it can be one of `BgWorkerStart_PostmasterStart` (start as soon as `postgres` itself has finished its own initialization; processes requesting this are not eligible for database connections), `BgWorkerStart_ConsistentState` (start as soon as a consistent state has been reached in a hot standby, allowing processes to connect to databases and run read-only queries), and `BgWorkerStart_RecoveryFinished` (start as soon as the system has entered normal read-write state). Note the last two values are equivalent in a server that's not a hot standby. Note that this setting only indicates when the processes are to be started; they do not stop when a different state is reached.

`bgw_restart_time` is the interval, in seconds, that `postgres` should wait before restarting the process in the event that it crashes. It can be any positive value, or `BGW_NEVER_RESTART`, indicating not to restart the process in case of a crash.

`bgw_library_name` is the name of a library in which the initial entry point for the background worker should be sought. The named library will be dynamically loaded by the worker process and `bgw_function_name` will be used to identify the function to be called. If calling a function in the core code, this must be set to `"postgres"`.

`bgw_function_name` is the name of the function to use as the initial entry point for the new background worker. If this function is in a dynamically loaded library, it must be marked `PGDLLEXPORT` (and not `static`).

`bgw_main_arg` is the `Datum` argument to the background worker main function. This main function should take a single argument of type `Datum` and return `void`. `bgw_main_arg` will be passed as the argument. In addition, the global variable `MyBgworkerEntry` points to a copy of the `BackgroundWorker` structure passed at registration time; the worker may find it helpful to examine this structure.

On Windows (and anywhere else where `EXEC_BACKEND` is defined) or in dynamic background workers it is not safe to pass a `Datum` by reference, only by value. If an argument is required, it is safest to pass an `int32` or other small value and use that as an index into an array allocated in shared memory. If a value like a `cstring` or `text` is passed then the pointer won't be valid from the new background worker process.

`bgw_extra` can contain extra data to be passed to the background worker. Unlike `bgw_main_arg`, this data is not passed as an argument to the worker's main function, but it can be accessed via `MyBgworkerEntry`, as discussed above.

`bgw_notify_pid` is the PID of a Postgres Pro backend process to which the postmaster should send `SIGUSR1` when the process is started or exits. It should be 0 for workers registered at postmaster startup time, or when the backend registering the worker does not wish to wait for the worker to start up. Otherwise, it should be initialized to `MyProcPid`.

Once running, the process can connect to a database by calling `BackgroundWorkerInitializeConnection(char *dbname, char *username, uint32 flags)` or `BackgroundWorkerInitializeConnectionByOid(Oid dboid, Oid useroid, uint32 flags)`. This allows the process to run transactions and queries using the SPI interface. If `dbname` is `NULL` or `dboid` is `InvalidOid`, the session is not connected to any particular database, but shared catalogs can be accessed. If `username` is `NULL` or `useroid` is `InvalidOid`, the process will run as the superuser created during `initdb`. If `BGWORKER_BYPASS_ALLOWCONN` is specified as `flags` it is possible to bypass the restriction to connect to databases not allowing user connections. If `BGWORKER_BYPASS_ROLELOGINCHECK` is specified as `flags` it is possible to bypass the login check for the role used to connect to databases. A background worker can only call one of these two functions, and only once. It is not possible to switch databases.

Signals are initially blocked when control reaches the background worker's main function, and must be unblocked by it; this is to allow the process to customize its signal handlers, if necessary. Signals can be unblocked in the new process by calling `BackgroundWorkerUnblockSignals` and blocked by calling `BackgroundWorkerBlockSignals`.

If `bgw_restart_time` for a background worker is configured as `BGW_NEVER_RESTART`, or if it exits with an exit code of 0 or is terminated by `TerminateBackgroundWorker`, it will be automatically unregis-

tered by the postmaster on exit. Otherwise, it will be restarted after the time period configured via `bgw_restart_time`, or immediately if the postmaster reinitializes the cluster due to a backend failure. Backends which need to suspend execution only temporarily should use an interruptible sleep rather than exiting; this can be achieved by calling `WaitLatch()`. Make sure the `WL_POSTMASTER_DEATH` flag is set when calling that function, and verify the return code for a prompt exit in the emergency case that `postgres` itself has terminated.

When a background worker is registered using the `RegisterDynamicBackgroundWorker` function, it is possible for the backend performing the registration to obtain information regarding the status of the worker. Backends wishing to do this should pass the address of a `BackgroundWorkerHandle *` as the second argument to `RegisterDynamicBackgroundWorker`. If the worker is successfully registered, this pointer will be initialized with an opaque handle that can subsequently be passed to `GetBackgroundWorkerPid(BackgroundWorkerHandle *, pid_t *)` or `TerminateBackgroundWorker(BackgroundWorkerHandle *)`. `GetBackgroundWorkerPid` can be used to poll the status of the worker: a return value of `BGWH_NOT_YET_STARTED` indicates that the worker has not yet been started by the postmaster; `BGWH_STOPPED` indicates that it has been started but is no longer running; and `BGWH_STARTED` indicates that it is currently running. In this last case, the PID will also be returned via the second argument. `TerminateBackgroundWorker` causes the postmaster to send `SIGTERM` to the worker if it is running, and to unregister it as soon as it is not.

In some cases, a process which registers a background worker may wish to wait for the worker to start up. This can be accomplished by initializing `bgw_notify_pid` to `MyProcPid` and then passing the `BackgroundWorkerHandle *` obtained at registration time to `WaitForBackgroundWorkerStartup(BackgroundWorkerHandle *handle, pid_t *)` function. This function will block until the postmaster has attempted to start the background worker, or until the postmaster dies. If the background worker is running, the return value will be `BGWH_STARTED`, and the PID will be written to the provided address. Otherwise, the return value will be `BGWH_STOPPED` or `BGWH_POSTMASTER_DIED`.

A process can also wait for a background worker to shut down, by using the `WaitForBackgroundWorkerShutdown(BackgroundWorkerHandle *handle)` function and passing the `BackgroundWorkerHandle *` obtained at registration. This function will block until the background worker exits, or postmaster dies. When the background worker exits, the return value is `BGWH_STOPPED`, if postmaster dies it will return `BGWH_POSTMASTER_DIED`.

Background workers can send asynchronous notification messages, either by using the `NOTIFY` command via SPI, or directly via `Async_Notify()`. Such notifications will be sent at transaction commit. Background workers should not register to receive asynchronous notifications with the `LISTEN` command, as there is no infrastructure for a worker to consume such notifications.

The maximum number of registered background workers is limited by [max_worker_processes](#).

Chapter 50. Logical Decoding

Postgres Pro provides infrastructure to stream the modifications performed via SQL to external consumers. This functionality can be used for a variety of purposes, including replication solutions and auditing.

Changes are sent out in streams identified by logical replication slots.

The format in which those changes are streamed is determined by the output plugin used. An example plugin is provided in the Postgres Pro distribution. Additional plugins can be written to extend the choice of available formats without modifying any core code. Every output plugin has access to each individual new row produced by `INSERT` and the new row version created by `UPDATE`. Availability of old row versions for `UPDATE` and `DELETE` depends on the configured replica identity (see [REPLICA IDENTITY](#)).

Changes can be consumed either using the streaming replication protocol (see [Section 56.4](#) and [Section 50.3](#)), or by calling functions via SQL (see [Section 50.4](#)). It is also possible to write additional methods of consuming the output of a replication slot without modifying core code (see [Section 50.7](#)).

50.1. Logical Decoding Examples

The following example demonstrates controlling logical decoding using the SQL interface.

Before you can use logical decoding, you must set `wal_level` to `logical` and `max_replication_slots` to at least 1. Then, you should connect to the target database (in the example below, `postgres`) as a superuser.

```
postgres=# -- Create a slot named 'regression_slot' using the output plugin
'test_decoding'
postgres=# SELECT * FROM pg_create_logical_replication_slot('regression_slot',
'test_decoding', false, true);
   slot_name   |    lsn
-----+-----
 regression_slot | 0/16B1970
(1 row)

postgres=# SELECT slot_name, plugin, slot_type, database, active, restart_lsn,
confirmed_flush_lsn FROM pg_replication_slots;
   slot_name   |   plugin   | slot_type | database | active | restart_lsn | confirmed_flush_lsn
-----+-----+-----+-----+-----+-----+-----
 regression_slot | test_decoding | logical   | postgres | f      | 0/16A4408   | 0/16A4440
(1 row)

postgres=# -- There are no changes to see yet
postgres=# SELECT * FROM pg_logical_slot_get_changes('regression_slot', NULL, NULL);
   lsn | xid | data
-----+-----+-----
(0 rows)

postgres=# CREATE TABLE data(id serial primary key, data text);
CREATE TABLE

postgres=# -- DDL isn't replicated, so all you'll see is the transaction
postgres=# SELECT * FROM pg_logical_slot_get_changes('regression_slot', NULL, NULL);
   lsn   | xid  | data
-----+-----+-----
 0/BA2DA58 | 10297 | BEGIN 10297
 0/BA5A5A0 | 10297 | COMMIT 10297
```

(2 rows)

```
postgres=# -- Once changes are read, they're consumed and not emitted
postgres=# -- in a subsequent call:
postgres=# SELECT * FROM pg_logical_slot_get_changes('regression_slot', NULL, NULL);
 lsn | xid | data
-----+-----+-----
(0 rows)
```

```
postgres=# BEGIN;
postgres=# INSERT INTO data(data) VALUES('1');
postgres=# INSERT INTO data(data) VALUES('2');
postgres=# COMMIT;
```

```
postgres=# SELECT * FROM pg_logical_slot_get_changes('regression_slot', NULL, NULL);
 lsn | xid | data
-----+-----+-----
0/BA5A688 | 10298 | BEGIN 10298
0/BA5A6F0 | 10298 | table public.data: INSERT: id[integer]:1 data[text]:'1'
0/BA5A7F8 | 10298 | table public.data: INSERT: id[integer]:2 data[text]:'2'
0/BA5A8A8 | 10298 | COMMIT 10298
(4 rows)
```

```
postgres=# INSERT INTO data(data) VALUES('3');
```

```
postgres=# -- You can also peek ahead in the change stream without consuming changes
postgres=# SELECT * FROM pg_logical_slot_peek_changes('regression_slot', NULL, NULL);
 lsn | xid | data
-----+-----+-----
0/BA5A8E0 | 10299 | BEGIN 10299
0/BA5A8E0 | 10299 | table public.data: INSERT: id[integer]:3 data[text]:'3'
0/BA5A990 | 10299 | COMMIT 10299
(3 rows)
```

```
postgres=# -- The next call to pg_logical_slot_peek_changes() returns the same changes
again
postgres=# SELECT * FROM pg_logical_slot_peek_changes('regression_slot', NULL, NULL);
 lsn | xid | data
-----+-----+-----
0/BA5A8E0 | 10299 | BEGIN 10299
0/BA5A8E0 | 10299 | table public.data: INSERT: id[integer]:3 data[text]:'3'
0/BA5A990 | 10299 | COMMIT 10299
(3 rows)
```

```
postgres=# -- options can be passed to output plugin, to influence the formatting
postgres=# SELECT * FROM pg_logical_slot_peek_changes('regression_slot', NULL, NULL,
 'include-timestamp', 'on');
 lsn | xid | data
-----+-----+-----
0/BA5A8E0 | 10299 | BEGIN 10299
0/BA5A8E0 | 10299 | table public.data: INSERT: id[integer]:3 data[text]:'3'
0/BA5A990 | 10299 | COMMIT 10299 (at 2017-05-10 12:07:21.272494-04)
(3 rows)
```

```
postgres=# -- Remember to destroy a slot you no longer need to stop it consuming
server resources:
postgres=# SELECT pg_drop_replication_slot('regression_slot');
pg_drop_replication_slot
```



```
-----
(1 row)
```

The following examples show how logical decoding is controlled over the streaming replication protocol, using the program `pg_recvlogical` included in the Postgres Pro distribution. This requires that client authentication is set up to allow replication connections (see [Section 26.2.5.1](#)) and that `max_wal_senders` is set sufficiently high to allow an additional connection. The second example shows how to stream two-phase transactions. Before you use two-phase commands, you must set `max_prepared_transactions` to at least 1.

Example 1:

```
$ pg_recvlogical -d postgres --slot=test --create-slot
$ pg_recvlogical -d postgres --slot=test --start -f -
Control+Z
$ psql -d postgres -c "INSERT INTO data(data) VALUES('4');"
$ fg
BEGIN 693
table public.data: INSERT: id[integer]:4 data[text]:'4'
COMMIT 693
Control+C
$ pg_recvlogical -d postgres --slot=test --drop-slot
```

Example 2:

```
$ pg_recvlogical -d postgres --slot=test --create-slot --two-phase
$ pg_recvlogical -d postgres --slot=test --start -f -
Control+Z
$ psql -d postgres -c "BEGIN;INSERT INTO data(data) VALUES('5');PREPARE TRANSACTION
'test';"
$ fg
BEGIN 694
table public.data: INSERT: id[integer]:5 data[text]:'5'
PREPARE TRANSACTION 'test', txid 694
Control+Z
$ psql -d postgres -c "COMMIT PREPARED 'test';"
$ fg
COMMIT PREPARED 'test', txid 694
Control+C
$ pg_recvlogical -d postgres --slot=test --drop-slot
```

The following example shows SQL interface that can be used to decode prepared transactions. Before you use two-phase commit commands, you must set `max_prepared_transactions` to at least 1. You must also have set the two-phase parameter as 'true' while creating the slot using `pg_create_logical_replication_slot`. Note that we will stream the entire transaction after the commit if it is not already decoded.

```
postgres=# BEGIN;
postgres=# INSERT INTO data(data) VALUES('5');
postgres=# PREPARE TRANSACTION 'test_prepared1';

postgres=# SELECT * FROM pg_logical_slot_get_changes('regression_slot', NULL, NULL);
   lsn    | xid | data
-----+-----+-----
 0/1689DC0 | 529 | BEGIN 529
 0/1689DC0 | 529 | table public.data: INSERT: id[integer]:3 data[text]:'5'
 0/1689FC0 | 529 | PREPARE TRANSACTION 'test_prepared1', txid 529
(3 rows)

postgres=# COMMIT PREPARED 'test_prepared1';
postgres=# select * from pg_logical_slot_get_changes('regression_slot', NULL, NULL);
```

```

lsn      | xid | data
-----+-----+-----
0/168A060 | 529 | COMMIT PREPARED 'test_prepared1', txid 529
(4 row)

postgres=#-- you can also rollback a prepared transaction
postgres=# BEGIN;
postgres=# INSERT INTO data(data) VALUES('6');
postgres=# PREPARE TRANSACTION 'test_prepared2';
postgres=# select * from pg_logical_slot_get_changes('regression_slot', NULL, NULL);
lsn      | xid | data
-----+-----+-----
0/168A180 | 530 | BEGIN 530
0/168A1E8 | 530 | table public.data: INSERT: id[integer]:4 data[text]:'6'
0/168A430 | 530 | PREPARE TRANSACTION 'test_prepared2', txid 530
(3 rows)

postgres=# ROLLBACK PREPARED 'test_prepared2';
postgres=# select * from pg_logical_slot_get_changes('regression_slot', NULL, NULL);
lsn      | xid | data
-----+-----+-----
0/168A4B8 | 530 | ROLLBACK PREPARED 'test_prepared2', txid 530
(1 row)

```

50.2. Logical Decoding Concepts

50.2.1. Logical Decoding

Logical decoding is the process of extracting all persistent changes to a database's tables into a coherent, easy to understand format which can be interpreted without detailed knowledge of the database's internal state.

In Postgres Pro, logical decoding is implemented by decoding the contents of the [write-ahead log](#), which describe changes on a storage level, into an application-specific form such as a stream of tuples or SQL statements.

50.2.2. Replication Slots

In the context of logical replication, a slot represents a stream of changes that can be replayed to a client in the order they were made on the origin server. Each slot streams a sequence of changes from a single database.

Note

Postgres Pro also has streaming replication slots (see [Section 26.2.5](#)), but they are used somewhat differently there.

A replication slot has an identifier that is unique across all databases in a Postgres Pro cluster. Slots persist independently of the connection using them and are crash-safe.

A logical slot will emit each change just once in normal operation. The current position of each slot is persisted only at checkpoint, so in the case of a crash the slot may return to an earlier LSN, which will then cause recent changes to be sent again when the server restarts. Logical decoding clients are responsible for avoiding ill effects from handling the same message more than once. Clients may wish to record the last LSN they saw when decoding and skip over any repeated data or (when using the replication protocol) request that decoding start from that LSN rather than letting the server determine the start point. The Replication Progress Tracking feature is designed for this purpose, refer to [replication origins](#).

Multiple independent slots may exist for a single database. Each slot has its own state, allowing different consumers to receive changes from different points in the database change stream. For most applications, a separate slot will be required for each consumer.

A logical replication slot knows nothing about the state of the receiver(s). It's even possible to have multiple different receivers using the same slot at different times; they'll just get the changes following on from when the last receiver stopped consuming them. Only one receiver may consume changes from a slot at any given time.

A logical replication slot can also be created on a hot standby. To prevent `VACUUM` from removing required rows from the system catalogs, `hot_standby_feedback` should be set on the standby. In spite of that, if any required rows get removed, the slot gets invalidated. It's highly recommended to use a physical slot between the primary and the standby. Otherwise, `hot_standby_feedback` will work but only while the connection is alive (for example a node restart would break it). Then, the primary may delete system catalog rows that could be needed by the logical decoding on the standby (as it does not know about the `catalog_xmin` on the standby). Existing logical slots on standby also get invalidated if `wal_level` on the primary is reduced to less than `logical`. This is done as soon as the standby detects such a change in the WAL stream. It means that, for walsenders that are lagging (if any), some WAL records up to the `wal_level` parameter change on the primary won't be decoded.

Creation of a logical slot requires information about all the currently running transactions. On the primary, this information is available directly, but on a standby, this information has to be obtained from primary. Thus, slot creation may need to wait for some activity to happen on the primary. If the primary is idle, creating a logical slot on standby may take noticeable time. This can be sped up by calling the `pg_log_standby_snapshot` function on the primary.

Caution

Replication slots persist across crashes and know nothing about the state of their consumer(s). They will prevent removal of required resources even when there is no connection using them. This consumes storage because neither required WAL nor required rows from the system catalogs can be removed by `VACUUM` as long as they are required by a replication slot. In extreme cases this could cause the database to shut down to prevent transaction ID wraparound (see [Section 24.1.5](#)). So if a slot is no longer required it should be dropped.

50.2.3. Replication Slot Synchronization

The logical replication slots on the primary can be synchronized to the hot standby by using the `failover` parameter of `pg_create_logical_replication_slot`, or by using the `failover` option of `CREATE SUBSCRIPTION` during slot creation, and then calling `pg_sync_replication_slots` on the standby. By setting `sync_replication_slots` on the standby, the failover slots can be synchronized periodically in the `slotsync` worker. For the synchronization to work, it is mandatory to have a physical replication slot between the primary and the standby (i.e., `primary_slot_name` should be configured on the standby), and `hot_standby_feedback` must be enabled on the standby. It is also necessary to specify a valid `db-name` in the `primary_conninfo`. It's highly recommended that the said physical replication slot is named in `synchronized_standby_slots` list on the primary, to prevent the subscriber from consuming changes faster than the hot standby. Even when correctly configured, some latency is expected when sending changes to logical subscribers due to the waiting on slots named in `synchronized_standby_slots`. When `synchronized_standby_slots` is utilized, the primary server will not completely shut down until the corresponding standbys, associated with the physical replication slots specified in `synchronized_standby_slots`, have confirmed receiving the WAL up to the latest flushed position on the primary server.

The ability to resume logical replication after failover depends upon the `pg_replication_slots.synced` value for the synchronized slots on the standby at the time of failover. Only persistent slots that have attained synced state as `true` on the standby before failover can be used for logical replication after failover. Temporary synced slots cannot be used for logical decoding, therefore logical replication for

those slots cannot be resumed. For example, if the synchronized slot could not become persistent on the standby due to a disabled subscription, then the subscription cannot be resumed after failover even when it is enabled.

To resume logical replication after failover from the synced logical slots, the subscription's 'conninfo' must be altered to point to the new primary server. This is done using `ALTER SUBSCRIPTION ... CONNECTION`. It is recommended that subscriptions are first disabled before promoting the standby and are re-enabled after altering the connection string.

Caution

There is a chance that the old primary is up again during the promotion and if subscriptions are not disabled, the logical subscribers may continue to receive data from the old primary server even after promotion until the connection string is altered. This might result in data inconsistency issues, preventing the logical subscribers from being able to continue replication from the new primary server.

50.2.4. Output Plugins

Output plugins transform the data from the write-ahead log's internal representation into the format the consumer of a replication slot desires.

50.2.5. Exported Snapshots

When a new replication slot is created using the streaming replication interface (see `CREATE_REPLICATION_SLOT`), a snapshot is exported (see [Section 9.28.5](#)), which will show exactly the state of the database after which all changes will be included in the change stream. This can be used to create a new replica by using `SET TRANSACTION SNAPSHOT` to read the state of the database at the moment the slot was created. This transaction can then be used to dump the database's state at that point in time, which afterwards can be updated using the slot's contents without losing any changes.

Creation of a snapshot is not always possible. In particular, it will fail when connected to a hot standby. Applications that do not require snapshot export may suppress it with the `NOEXPORT_SNAPSHOT` option.

50.3. Streaming Replication Protocol Interface

The commands

- `CREATE_REPLICATION_SLOT slot_name LOGICAL output_plugin`
- `DROP_REPLICATION_SLOT slot_name [WAIT]`
- `START_REPLICATION SLOT slot_name LOGICAL ...`

are used to create, drop, and stream changes from a replication slot, respectively. These commands are only available over a replication connection; they cannot be used via SQL. See [Section 56.4](#) for details on these commands.

The command `pg_recvlogical` can be used to control logical decoding over a streaming replication connection. (It uses these commands internally.)

50.4. Logical Decoding SQL Interface

See [Section 9.28.6](#) for detailed documentation on the SQL-level API for interacting with logical decoding.

Synchronous replication (see [Section 26.2.8](#)) is only supported on replication slots used over the streaming replication interface. The function interface and additional, non-core interfaces do not support synchronous replication.

50.5. System Catalogs Related to Logical Decoding

The `pg_replication_slots` view and the `pg_stat_replication` view provide information about the current state of replication slots and streaming replication connections respectively. These views apply to both physical and logical replication. The `pg_stat_replication_slots` view provides statistics information about the logical replication slots.

50.6. Logical Decoding Output Plugins

50.6.1. Initialization Function

An output plugin is loaded by dynamically loading a shared library with the output plugin's name as the library base name. The normal library search path is used to locate the library. To provide the required output plugin callbacks and to indicate that the library is actually an output plugin it needs to provide a function named `_PG_output_plugin_init`. This function is passed a struct that needs to be filled with the callback function pointers for individual actions.

```
typedef struct OutputPluginCallbacks
{
    LogicalDecodeStartupCB startup_cb;
    LogicalDecodeBeginCB begin_cb;
    LogicalDecodeChangeCB change_cb;
    LogicalDecodeTruncateCB truncate_cb;
    LogicalDecodeCommitCB commit_cb;
    LogicalDecodeMessageCB message_cb;
    LogicalDecodeFilterByOriginCB filter_by_origin_cb;
    LogicalDecodeShutdownCB shutdown_cb;
    LogicalDecodeFilterPrepareCB filter_prepare_cb;
    LogicalDecodeBeginPrepareCB begin_prepare_cb;
    LogicalDecodePrepareCB prepare_cb;
    LogicalDecodeCommitPreparedCB commit_prepared_cb;
    LogicalDecodeRollbackPreparedCB rollback_prepared_cb;
    LogicalDecodeStreamStartCB stream_start_cb;
    LogicalDecodeStreamStopCB stream_stop_cb;
    LogicalDecodeStreamAbortCB stream_abort_cb;
    LogicalDecodeStreamPrepareCB stream_prepare_cb;
    LogicalDecodeStreamCommitCB stream_commit_cb;
    LogicalDecodeStreamChangeCB stream_change_cb;
    LogicalDecodeStreamMessageCB stream_message_cb;
    LogicalDecodeStreamTruncateCB stream_truncate_cb;
} OutputPluginCallbacks;
```

```
typedef void (*LogicalOutputPluginInit) (struct OutputPluginCallbacks *cb);
```

The `begin_cb`, `change_cb` and `commit_cb` callbacks are required, while `startup_cb`, `truncate_cb`, `message_cb`, `filter_by_origin_cb`, and `shutdown_cb` are optional. If `truncate_cb` is not set but a `TRUNCATE` is to be decoded, the action will be ignored.

An output plugin may also define functions to support streaming of large, in-progress transactions. The `stream_start_cb`, `stream_stop_cb`, `stream_abort_cb`, `stream_commit_cb`, and `stream_change_cb` are required, while `stream_message_cb` and `stream_truncate_cb` are optional. The `stream_prepare_cb` is also required if the output plugin also support two-phase commits.

An output plugin may also define functions to support two-phase commits, which allows actions to be decoded on the `PREPARE TRANSACTION`. The `begin_prepare_cb`, `prepare_cb`, `commit_prepared_cb` and `rollback_prepared_cb` callbacks are required, while `filter_prepare_cb` is optional. The `stream_prepare_cb` is also required if the output plugin also supports the streaming of large in-progress transactions.

50.6.2. Capabilities

To decode, format and output changes, output plugins can use most of the backend's normal infrastructure, including calling output functions. Read only access to relations is permitted as long as only relations are accessed that either have been created by `initdb` in the `pg_catalog` schema, or have been marked as user provided catalog tables using

```
ALTER TABLE user_catalog_table SET (user_catalog_table = true);
CREATE TABLE another_catalog_table(data text) WITH (user_catalog_table = true);
```

Note that access to user catalog tables or regular system catalog tables in the output plugins has to be done via the `systable_*` scan APIs only. Access via the `heap_*` scan APIs will error out. Additionally, any actions leading to transaction ID assignment are prohibited. That, among others, includes writing to tables, performing DDL changes, and calling `pg_current_xact_id()`.

50.6.3. Output Modes

Output plugin callbacks can pass data to the consumer in nearly arbitrary formats. For some use cases, like viewing the changes via SQL, returning data in a data type that can contain arbitrary data (e.g., `bytea`) is cumbersome. If the output plugin only outputs textual data in the server's encoding, it can declare that by setting `OutputPluginOptions.output_type` to `OUTPUT_PLUGIN_TEXTUAL_OUTPUT` instead of `OUTPUT_PLUGIN_BINARY_OUTPUT` in the [startup callback](#). In that case, all the data has to be in the server's encoding so that a `text` datum can contain it. This is checked in assertion-enabled builds.

50.6.4. Output Plugin Callbacks

An output plugin gets notified about changes that are happening via various callbacks it needs to provide.

Concurrent transactions are decoded in commit order, and only changes belonging to a specific transaction are decoded between the `begin` and `commit` callbacks. Transactions that were rolled back explicitly or implicitly never get decoded. Successful savepoints are folded into the transaction containing them in the order they were executed within that transaction. A transaction that is prepared for a two-phase commit using `PREPARE TRANSACTION` will also be decoded if the output plugin callbacks needed for decoding them are provided. It is possible that the current prepared transaction which is being decoded is aborted concurrently via a `ROLLBACK PREPARED` command. In that case, the logical decoding of this transaction will be aborted too. All the changes of such a transaction are skipped once the abort is detected and the `prepare_cb` callback is invoked. Thus even in case of a concurrent abort, enough information is provided to the output plugin for it to properly deal with `ROLLBACK PREPARED` once that is decoded.

Note

Only transactions that have already safely been flushed to disk will be decoded. That can lead to a `COMMIT` not immediately being decoded in a directly following `pg_logical_slot_get_changes()` when `synchronous_commit` is set to `off`.

50.6.4.1. Startup Callback

The optional `startup_cb` callback is called whenever a replication slot is created or asked to stream changes, independent of the number of changes that are ready to be put out.

```
typedef void (*LogicalDecodeStartupCB) (struct LogicalDecodingContext *ctx,
                                         OutputPluginOptions *options,
                                         bool is_init);
```

The `is_init` parameter will be true when the replication slot is being created and false otherwise. `options` points to a struct of options that output plugins can set:

```
typedef struct OutputPluginOptions
{
    OutputPluginOutputType output_type;
```

```
bool        receive_rewrites;
} OutputPluginOptions;
```

`output_type` has to either be set to `OUTPUT_PLUGIN_TEXTUAL_OUTPUT` or `OUTPUT_PLUGIN_BINARY_OUTPUT`. See also [Section 50.6.3](#). If `receive_rewrites` is true, the output plugin will also be called for changes made by heap rewrites during certain DDL operations. These are of interest to plugins that handle DDL replication, but they require special handling.

The startup callback should validate the options present in `ctx->output_plugin_options`. If the output plugin needs to have a state, it can use `ctx->output_plugin_private` to store it.

50.6.4.2. Shutdown Callback

The optional `shutdown_cb` callback is called whenever a formerly active replication slot is not used anymore and can be used to deallocate resources private to the output plugin. The slot isn't necessarily being dropped, streaming is just being stopped.

```
typedef void (*LogicalDecodeShutdownCB) (struct LogicalDecodingContext *ctx);
```

50.6.4.3. Transaction Begin Callback

The required `begin_cb` callback is called whenever a start of a committed transaction has been decoded. Aborted transactions and their contents never get decoded.

```
typedef void (*LogicalDecodeBeginCB) (struct LogicalDecodingContext *ctx,
                                     ReorderBufferTXN *txn);
```

The `txn` parameter contains meta information about the transaction, like the time stamp at which it has been committed and its `XID`.

50.6.4.4. Transaction End Callback

The required `commit_cb` callback is called whenever a transaction commit has been decoded. The `change_cb` callbacks for all modified rows will have been called before this, if there have been any modified rows.

```
typedef void (*LogicalDecodeCommitCB) (struct LogicalDecodingContext *ctx,
                                       ReorderBufferTXN *txn,
                                       XLogRecPtr commit_lsn);
```

50.6.4.5. Change Callback

The required `change_cb` callback is called for every individual row modification inside a transaction, may it be an `INSERT`, `UPDATE`, or `DELETE`. Even if the original command modified several rows at once the callback will be called individually for each row. The `change_cb` callback may access system or user catalog tables to aid in the process of outputting the row modification details. In case of decoding a prepared (but yet uncommitted) transaction or decoding of an uncommitted transaction, this change callback might also error out due to simultaneous rollback of this very same transaction. In that case, the logical decoding of this aborted transaction is stopped gracefully.

```
typedef void (*LogicalDecodeChangeCB) (struct LogicalDecodingContext *ctx,
                                       ReorderBufferTXN *txn,
                                       Relation relation,
                                       ReorderBufferChange *change);
```

The `ctx` and `txn` parameters have the same contents as for the `begin_cb` and `commit_cb` callbacks, but additionally the relation descriptor `relation` points to the relation the row belongs to and a struct `change` describing the row modification are passed in.

Note

Only changes in user defined tables that are not unlogged (see [UNLOGGED](#)) and not temporary (see [TEMPORARY](#) or [TEMP](#)) can be extracted using logical decoding.

50.6.4.6. Truncate Callback

The optional `truncate_cb` callback is called for a `TRUNCATE` command.

```
typedef void (*LogicalDecodeTruncateCB) (struct LogicalDecodingContext *ctx,
                                         ReorderBufferTXN *txn,
                                         int nrelations,
                                         Relation relations[],
                                         ReorderBufferChange *change);
```

The parameters are analogous to the `change_cb` callback. However, because `TRUNCATE` actions on tables connected by foreign keys need to be executed together, this callback receives an array of relations instead of just a single one. See the description of the [TRUNCATE](#) statement for details.

50.6.4.7. Origin Filter Callback

The optional `filter_by_origin_cb` callback is called to determine whether data that has been replayed from `origin_id` is of interest to the output plugin.

```
typedef bool (*LogicalDecodeFilterByOriginCB) (struct LogicalDecodingContext *ctx,
                                              RepOriginId origin_id);
```

The `ctx` parameter has the same contents as for the other callbacks. No information but the origin is available. To signal that changes originating on the passed in node are irrelevant, return true, causing them to be filtered away; false otherwise. The other callbacks will not be called for transactions and changes that have been filtered away.

This is useful when implementing cascading or multidirectional replication solutions. Filtering by the origin allows to prevent replicating the same changes back and forth in such setups. While transactions and changes also carry information about the origin, filtering via this callback is noticeably more efficient.

50.6.4.8. Generic Message Callback

The optional `message_cb` callback is called whenever a logical decoding message has been decoded.

```
typedef void (*LogicalDecodeMessageCB) (struct LogicalDecodingContext *ctx,
                                         ReorderBufferTXN *txn,
                                         XLogRecPtr message_lsn,
                                         bool transactional,
                                         const char *prefix,
                                         Size message_size,
                                         const char *message);
```

The `txn` parameter contains meta information about the transaction, like the time stamp at which it has been committed and its XID. Note however that it can be NULL when the message is non-transactional and the XID was not assigned yet in the transaction which logged the message. The `lsn` has WAL location of the message. The `transactional` says if the message was sent as transactional or not. Similar to the change callback, in case of decoding a prepared (but yet uncommitted) transaction or decoding of an uncommitted transaction, this message callback might also error out due to simultaneous rollback of this very same transaction. In that case, the logical decoding of this aborted transaction is stopped gracefully. The `prefix` is arbitrary null-terminated prefix which can be used for identifying interesting messages for the current plugin. And finally the `message` parameter holds the actual message of `message_size` size.

Extra care should be taken to ensure that the prefix the output plugin considers interesting is unique. Using name of the extension or the output plugin itself is often a good choice.

50.6.4.9. Prepare Filter Callback

The optional `filter_prepare_cb` callback is called to determine whether data that is part of the current two-phase commit transaction should be considered for decoding at this prepare stage or later as a regular one-phase transaction at `COMMIT PREPARED` time. To signal that decoding should be skipped,

`return true; false` otherwise. When the callback is not defined, `false` is assumed (i.e. no filtering, all transactions using two-phase commit are decoded in two phases as well).

```
typedef bool (*LogicalDecodeFilterPrepareCB) (struct LogicalDecodingContext *ctx,
                                             TransactionId xid,
                                             const char *gid);
```

The `ctx` parameter has the same contents as for the other callbacks. The parameters `xid` and `gid` provide two different ways to identify the transaction. The later `COMMIT PREPARED` or `ROLLBACK PREPARED` carries both identifiers, providing an output plugin the choice of what to use.

The callback may be invoked multiple times per transaction to decode and must provide the same static answer for a given pair of `xid` and `gid` every time it is called.

50.6.4.10. Transaction Begin Prepare Callback

The required `begin_prepare_cb` callback is called whenever the start of a prepared transaction has been decoded. The `gid` field, which is part of the `txn` parameter, can be used in this callback to check if the plugin has already received this `PREPARE` in which case it can either error out or skip the remaining changes of the transaction.

```
typedef void (*LogicalDecodeBeginPrepareCB) (struct LogicalDecodingContext *ctx,
                                             ReorderBufferTXN *txn);
```

50.6.4.11. Transaction Prepare Callback

The required `prepare_cb` callback is called whenever a transaction which is prepared for two-phase commit has been decoded. The `change_cb` callback for all modified rows will have been called before this, if there have been any modified rows. The `gid` field, which is part of the `txn` parameter, can be used in this callback.

```
typedef void (*LogicalDecodePrepareCB) (struct LogicalDecodingContext *ctx,
                                        ReorderBufferTXN *txn,
                                        XLogRecPtr prepare_lsn);
```

50.6.4.12. Transaction Commit Prepared Callback

The required `commit_prepared_cb` callback is called whenever a transaction `COMMIT PREPARED` has been decoded. The `gid` field, which is part of the `txn` parameter, can be used in this callback.

```
typedef void (*LogicalDecodeCommitPreparedCB) (struct LogicalDecodingContext *ctx,
                                                ReorderBufferTXN *txn,
                                                XLogRecPtr commit_lsn);
```

50.6.4.13. Transaction Rollback Prepared Callback

The required `rollback_prepared_cb` callback is called whenever a transaction `ROLLBACK PREPARED` has been decoded. The `gid` field, which is part of the `txn` parameter, can be used in this callback. The parameters `prepare_end_lsn` and `prepare_time` can be used to check if the plugin has received this `PREPARE TRANSACTION` in which case it can apply the rollback, otherwise, it can skip the rollback operation. The `gid` alone is not sufficient because the downstream node can have a prepared transaction with same identifier.

```
typedef void (*LogicalDecodeRollbackPreparedCB) (struct LogicalDecodingContext *ctx,
                                                  ReorderBufferTXN *txn,
                                                  XLogRecPtr prepare_end_lsn,
                                                  TimestampTz prepare_time);
```

50.6.4.14. Stream Start Callback

The required `stream_start_cb` callback is called when opening a block of streamed changes from an in-progress transaction.

```
typedef void (*LogicalDecodeStreamStartCB) (struct LogicalDecodingContext *ctx,
                                             ReorderBufferTXN *txn);
```

50.6.4.15. Stream Stop Callback

The required `stream_stop_cb` callback is called when closing a block of streamed changes from an in-progress transaction.

```
typedef void (*LogicalDecodeStreamStopCB) (struct LogicalDecodingContext *ctx,
                                           ReorderBufferTXN *txn);
```

50.6.4.16. Stream Abort Callback

The required `stream_abort_cb` callback is called to abort a previously streamed transaction.

```
typedef void (*LogicalDecodeStreamAbortCB) (struct LogicalDecodingContext *ctx,
                                           ReorderBufferTXN *txn,
                                           XLogRecPtr abort_lsn);
```

50.6.4.17. Stream Prepare Callback

The `stream_prepare_cb` callback is called to prepare a previously streamed transaction as part of a two-phase commit. This callback is required when the output plugin supports both the streaming of large in-progress transactions and two-phase commits.

```
typedef void (*LogicalDecodeStreamPrepareCB) (struct LogicalDecodingContext *ctx,
                                              ReorderBufferTXN *txn,
                                              XLogRecPtr prepare_lsn);
```

50.6.4.18. Stream Commit Callback

The required `stream_commit_cb` callback is called to commit a previously streamed transaction.

```
typedef void (*LogicalDecodeStreamCommitCB) (struct LogicalDecodingContext *ctx,
                                             ReorderBufferTXN *txn,
                                             XLogRecPtr commit_lsn);
```

50.6.4.19. Stream Change Callback

The required `stream_change_cb` callback is called when sending a change in a block of streamed changes (demarcated by `stream_start_cb` and `stream_stop_cb` calls). The actual changes are not displayed as the transaction can abort at a later point in time and we don't decode changes for aborted transactions.

```
typedef void (*LogicalDecodeStreamChangeCB) (struct LogicalDecodingContext *ctx,
                                             ReorderBufferTXN *txn,
                                             Relation relation,
                                             ReorderBufferChange *change);
```

50.6.4.20. Stream Message Callback

The optional `stream_message_cb` callback is called when sending a generic message in a block of streamed changes (demarcated by `stream_start_cb` and `stream_stop_cb` calls). The message contents for transactional messages are not displayed as the transaction can abort at a later point in time and we don't decode changes for aborted transactions.

```
typedef void (*LogicalDecodeStreamMessageCB) (struct LogicalDecodingContext *ctx,
                                              ReorderBufferTXN *txn,
                                              XLogRecPtr message_lsn,
                                              bool transactional,
                                              const char *prefix,
                                              Size message_size,
                                              const char *message);
```

50.6.4.21. Stream Truncate Callback

The optional `stream_truncate_cb` callback is called for a `TRUNCATE` command in a block of streamed changes (demarcated by `stream_start_cb` and `stream_stop_cb` calls).

```
typedef void (*LogicalDecodeStreamTruncateCB) (struct LogicalDecodingContext *ctx,
                                              ReorderBufferTXN *txn,
                                              int nrelations,
                                              Relation relations[],
                                              ReorderBufferChange *change);
```

The parameters are analogous to the `stream_change_cb` callback. However, because `TRUNCATE` actions on tables connected by foreign keys need to be executed together, this callback receives an array of relations instead of just a single one. See the description of the [TRUNCATE](#) statement for details.

50.6.5. Functions for Producing Output

To actually produce output, output plugins can write data to the `StringInfo` output buffer in `ctx->out` when inside the `begin_cb`, `commit_cb`, or `change_cb` callbacks. Before writing to the output buffer, `OutputPluginPrepareWrite(ctx, last_write)` has to be called, and after finishing writing to the buffer, `OutputPluginWrite(ctx, last_write)` has to be called to perform the write. The `last_write` indicates whether a particular write was the callback's last write.

The following example shows how to output data to the consumer of an output plugin:

```
OutputPluginPrepareWrite(ctx, true);
appendStringInfo(ctx->out, "BEGIN %u", txn->xid);
OutputPluginWrite(ctx, true);
```

50.7. Logical Decoding Output Writers

It is possible to add more output methods for logical decoding. Essentially, three functions need to be provided: one to read WAL, one to prepare writing output, and one to write the output (see [Section 50.6.5](#)).

50.8. Synchronous Replication Support for Logical Decoding

50.8.1. Overview

Logical decoding can be used to build [synchronous replication](#) solutions with the same user interface as synchronous replication for [streaming replication](#). To do this, the streaming replication interface (see [Section 50.3](#)) must be used to stream out data. Clients have to send `Standby status update (F)` (see [Section 56.4](#)) messages, just like streaming replication clients do.

Note

A synchronous replica receiving changes via logical decoding will work in the scope of a single database. Since, in contrast to that, `synchronous_standby_names` currently is server wide, this means this technique will not work properly if more than one database is actively used.

50.8.2. Caveats

In synchronous replication setup, a deadlock can happen, if the transaction has locked [user] catalog tables exclusively. See [Section 50.6.2](#) for information on user catalog tables. This is because logical decoding of transactions can lock catalog tables to access them. To avoid this users must refrain from taking an exclusive lock on [user] catalog tables. This can happen in the following ways:

- Issuing an explicit `LOCK` on `pg_class` in a transaction.
- Perform `CLUSTER` on `pg_class` in a transaction.
- `PREPARE TRANSACTION` after `LOCK` command on `pg_class` and allow logical decoding of two-phase transactions.

- `PREPARE TRANSACTION` after `CLUSTER` command on `pg_trigger` and allow logical decoding of two-phase transactions. This will lead to deadlock only when published table have a trigger.
- Executing `TRUNCATE` on [user] catalog table in a transaction.

Note that these commands can cause deadlocks not only for the system catalog tables listed above but for other catalog tables.

50.9. Streaming of Large Transactions for Logical Decoding

The basic output plugin callbacks (e.g., `begin_cb`, `change_cb`, `commit_cb` and `message_cb`) are only invoked when the transaction actually commits. The changes are still decoded from the transaction log, but are only passed to the output plugin at commit (and discarded if the transaction aborts).

This means that while the decoding happens incrementally, and may spill to disk to keep memory usage under control, all the decoded changes have to be transmitted when the transaction finally commits (or more precisely, when the commit is decoded from the transaction log). Depending on the size of the transaction and network bandwidth, the transfer time may significantly increase the apply lag.

To reduce the apply lag caused by large transactions, an output plugin may provide additional callback to support incremental streaming of in-progress transactions. There are multiple required streaming callbacks (`stream_start_cb`, `stream_stop_cb`, `stream_abort_cb`, `stream_commit_cb` and `stream_change_cb`) and two optional callbacks (`stream_message_cb` and `stream_truncate_cb`). Also, if streaming of two-phase commands is to be supported, then additional callbacks must be provided. (See [Section 50.10](#) for details).

When streaming an in-progress transaction, the changes (and messages) are streamed in blocks demarcated by `stream_start_cb` and `stream_stop_cb` callbacks. Once all the decoded changes are transmitted, the transaction can be committed using the `stream_commit_cb` callback (or possibly aborted using the `stream_abort_cb` callback). If two-phase commits are supported, the transaction can be prepared using the `stream_prepare_cb` callback, `COMMIT PREPARED` using the `commit_prepared_cb` callback or aborted using the `rollback_prepared_cb`.

One example sequence of streaming callback calls for one transaction may look like this:

```
stream_start_cb(...);    <-- start of first block of changes
  stream_change_cb(...);
  stream_change_cb(...);
  stream_message_cb(...);
  stream_change_cb(...);
  ...
  stream_change_cb(...);
stream_stop_cb(...);     <-- end of first block of changes

stream_start_cb(...);    <-- start of second block of changes
  stream_change_cb(...);
  stream_change_cb(...);
  stream_change_cb(...);
  ...
  stream_message_cb(...);
  stream_change_cb(...);
stream_stop_cb(...);     <-- end of second block of changes

[a. when using normal commit]
stream_commit_cb(...);    <-- commit of the streamed transaction

[b. when using two-phase commit]
```

```
stream_prepare_cb(...);    <-- prepare the streamed transaction
commit_prepared_cb(...);  <-- commit of the prepared transaction
```

The actual sequence of callback calls may be more complicated, of course. There may be blocks for multiple streamed transactions, some of the transactions may get aborted, etc.

Similar to spill-to-disk behavior, streaming is triggered when the total amount of changes decoded from the WAL (for all in-progress transactions) exceeds the limit defined by `logical_decoding_work_mem` setting. At that point, the largest top-level transaction (measured by the amount of memory currently used for decoded changes) is selected and streamed. However, in some cases we still have to spill to disk even if streaming is enabled because we exceed the memory threshold but still have not decoded the complete tuple e.g., only decoded toast table insert but not the main table insert.

Even when streaming large transactions, the changes are still applied in commit order, preserving the same guarantees as the non-streaming mode.

50.10. Two-phase Commit Support for Logical Decoding

With the basic output plugin callbacks (eg., `begin_cb`, `change_cb`, `commit_cb` and `message_cb`) two-phase commit commands like `PREPARE TRANSACTION`, `COMMIT PREPARED` and `ROLLBACK PREPARED` are not decoded. While the `PREPARE TRANSACTION` is ignored, `COMMIT PREPARED` is decoded as a `COMMIT` and `ROLLBACK PREPARED` is decoded as a `ROLLBACK`.

To support the streaming of two-phase commands, an output plugin needs to provide additional callbacks. There are multiple two-phase commit callbacks that are required, (`begin_prepare_cb`, `prepare_cb`, `commit_prepared_cb`, `rollback_prepared_cb` and `stream_prepare_cb`) and an optional callback (`filter_prepare_cb`).

If the output plugin callbacks for decoding two-phase commit commands are provided, then on `PREPARE TRANSACTION`, the changes of that transaction are decoded, passed to the output plugin, and the `prepare_cb` callback is invoked. This differs from the basic decoding setup where changes are only passed to the output plugin when a transaction is committed. The start of a prepared transaction is indicated by the `begin_prepare_cb` callback.

When a prepared transaction is rolled back using the `ROLLBACK PREPARED`, then the `rollback_prepared_cb` callback is invoked and when the prepared transaction is committed using `COMMIT PREPARED`, then the `commit_prepared_cb` callback is invoked.

Optionally the output plugin can define filtering rules via `filter_prepare_cb` to decode only specific transaction in two phases. This can be achieved by pattern matching on the `gid` or via lookups using the `xid`.

The users that want to decode prepared transactions need to be careful about below mentioned points:

- If the prepared transaction has locked [user] catalog tables exclusively then decoding prepare can block till the main transaction is committed.
- The logical replication solution that builds distributed two phase commit using this feature can deadlock if the prepared transaction has locked [user] catalog tables exclusively. To avoid this users must refrain from having locks on catalog tables (e.g. explicit `LOCK` command) in such transactions. See [Section 50.8.2](#) for the details.

Chapter 51. Replication Progress Tracking

Replication origins are intended to make it easier to implement logical replication solutions on top of [logical decoding](#). They provide a solution to two common problems:

- How to safely keep track of replication progress
- How to change replication behavior based on the origin of a row; for example, to prevent loops in bi-directional replication setups

Replication origins have just two properties, a name and an ID. The name, which is what should be used to refer to the origin across systems, is free-form `text`. It should be used in a way that makes conflicts between replication origins created by different replication solutions unlikely; e.g., by prefixing the replication solution's name to it. The ID is used only to avoid having to store the long version in situations where space efficiency is important. It should never be shared across systems.

Replication origins can be created using the function `pg_replication_origin_create()`; dropped using `pg_replication_origin_drop()`; and seen in the `pg_replication_origin` system catalog.

One nontrivial part of building a replication solution is to keep track of replay progress in a safe manner. When the applying process, or the whole cluster, dies, it needs to be possible to find out up to where data has successfully been replicated. Naive solutions to this, such as updating a row in a table for every replayed transaction, have problems like run-time overhead and database bloat.

Using the replication origin infrastructure a session can be marked as replaying from a remote node (using the `pg_replication_origin_session_setup()` function). Additionally the LSN and commit time stamp of every source transaction can be configured on a per transaction basis using `pg_replication_origin_xact_setup()`. If that's done replication progress will persist in a crash safe manner. Replay progress for all replication origins can be seen in the `pg_replication_origin_status` view. An individual origin's progress, e.g., when resuming replication, can be acquired using `pg_replication_origin_progress()` for any origin or `pg_replication_origin_session_progress()` for the origin configured in the current session.

In replication topologies more complex than replication from exactly one system to one other system, another problem can be that it is hard to avoid replicating replayed rows again. That can lead both to cycles in the replication and inefficiencies. Replication origins provide an optional mechanism to recognize and prevent that. When configured using the functions referenced in the previous paragraph, every change and transaction passed to output plugin callbacks (see [Section 50.6](#)) generated by the session is tagged with the replication origin of the generating session. This allows treating them differently in the output plugin, e.g., ignoring all but locally-originating rows. Additionally the `filter_by_origin_cb` callback can be used to filter the logical decoding change stream based on the source. While less flexible, filtering via that callback is considerably more efficient than doing it in the output plugin.

Chapter 52. Archive Modules

Postgres Pro provides infrastructure to create custom modules for continuous archiving (see [Section 25.1](#)). While archiving via a shell command (i.e., [archive_command](#)) is much simpler, a custom archive module will often be considerably more robust and performant.

When a custom [archive_library](#) is configured, Postgres Pro will submit completed WAL files to the module, and the server will avoid recycling or removing these WAL files until the module indicates that the files were successfully archived. It is ultimately up to the module to decide what to do with each WAL file, but many recommendations are listed at [Section 25.1.1](#).

Archiving modules must at least consist of an initialization function (see [Section 52.1](#)) and the required callbacks (see [Section 52.2](#)). However, archive modules are also permitted to do much more (e.g., declare GUCs and register background workers).

The `contrib/basic_archive` module contains a working example, which demonstrates some useful techniques.

52.1. Initialization Functions

An archive library is loaded by dynamically loading a shared library with the [archive_library](#)'s name as the library base name. The normal library search path is used to locate the library. To provide the required archive module callbacks and to indicate that the library is actually an archive module, it needs to provide a function named `_PG_archive_module_init`. The result of the function must be a pointer to a struct of type `ArchiveModuleCallbacks`, which contains everything that the core code needs to know to make use of the archive module. The return value needs to be of server lifetime, which is typically achieved by defining it as a `static const` variable in global scope.

```
typedef struct ArchiveModuleCallbacks
{
    ArchiveStartupCB startup_cb;
    ArchiveCheckConfiguredCB check_configured_cb;
    ArchiveFileCB archive_file_cb;
    ArchiveShutdownCB shutdown_cb;
} ArchiveModuleCallbacks;
typedef const ArchiveModuleCallbacks *(*ArchiveModuleInit) (void);
```

Only the `archive_file_cb` callback is required. The others are optional.

52.2. Archive Module Callbacks

The archive callbacks define the actual archiving behavior of the module. The server will call them as required to process each individual WAL file.

52.2.1. Startup Callback

The `startup_cb` callback is called shortly after the module is loaded. This callback can be used to perform any additional initialization required. If the archive module has any state, it can use `state->private_data` to store it.

```
typedef void (*ArchiveStartupCB) (ArchiveModuleState *state);
```

52.2.2. Check Callback

The `check_configured_cb` callback is called to determine whether the module is fully configured and ready to accept WAL files (e.g., its configuration parameters are set to valid values). If no `check_configured_cb` is defined, the server always assumes the module is configured.

```
typedef bool (*ArchiveCheckConfiguredCB) (ArchiveModuleState *state);
```

If `true` is returned, the server will proceed with archiving the file by calling the `archive_file_cb` callback. If `false` is returned, archiving will not proceed, and the archiver will emit the following message to the server log:

```
WARNING: archive_mode enabled, yet archiving is not configured
```

In the latter case, the server will periodically call this function, and archiving will proceed only when it returns `true`.

Note

When returning `false`, it may be useful to append some additional information to the generic warning message. To do that, provide a message to the `arch_module_check_errdetail` macro before returning `false`. Like `errdetail()`, this macro accepts a format string followed by an optional list of arguments. The resulting string will be emitted as the `DETAIL` line of the warning message.

52.2.3. Archive Callback

The `archive_file_cb` callback is called to archive a single WAL file.

```
typedef bool (*ArchiveFileCB) (ArchiveModuleState *state, const char *file, const char *path);
```

If `true` is returned, the server proceeds as if the file was successfully archived, which may include recycling or removing the original WAL file. If `false` is returned or an error is thrown, the server will keep the original WAL file and retry archiving later. `file` will contain just the file name of the WAL file to archive, while `path` contains the full path of the WAL file (including the file name).

Note

The `archive_file_cb` callback is called in a short-lived memory context that will be reset between invocations. If you need longer-lived storage, create a memory context in the module's `startup_cb` callback.

52.2.4. Shutdown Callback

The `shutdown_cb` callback is called when the archiver process exits (e.g., after an error) or the value of [archive_library](#) changes. If no `shutdown_cb` is defined, no special action is taken in these situations. If the archive module has any state, this callback should free it to avoid leaks.

```
typedef void (*ArchiveShutdownCB) (ArchiveModuleState *state);
```

Part VI. Reference

The entries in this Reference are meant to provide in reasonable length an authoritative, complete, and formal summary about their respective subjects. More information about the use of Postgres Pro, in narrative, tutorial, or example form, can be found in other parts of this book. See the cross-references listed on each reference page.

The reference entries are also available as traditional “man” pages.

SQL Commands

This part contains reference information for the SQL commands supported by Postgres Pro. By “SQL” the language in general is meant; information about the standards conformance and compatibility of each command can be found on the respective reference page.

ABORT

ABORT — abort the current transaction

Synopsis

```
ABORT [ AUTONOMOUS ] [ WORK | TRANSACTION ] [ AND [ NO ] CHAIN ]
```

Description

ABORT rolls back the current transaction and causes all the updates made by the transaction to be discarded. This command is identical in behavior to the standard SQL command [ROLLBACK](#), and is present only for historical reasons.

Parameters

AUTONOMOUS

Optional key word that can be used when aborting an autonomous transaction. For details on autonomous transactions, see [Chapter 16](#).

WORK

TRANSACTION

Optional key words. They have no effect.

AND CHAIN

If AND CHAIN is specified, a new transaction is immediately started with the same transaction characteristics (see [SET TRANSACTION](#)) as the just finished one. Otherwise, no new transaction is started.

Notes

Use [COMMIT](#) to successfully terminate a transaction.

Issuing ABORT outside of a transaction block emits a warning and otherwise has no effect.

Examples

To abort all changes:

```
ABORT;
```

Compatibility

This command is a Postgres Pro extension present for historical reasons. ROLLBACK is the equivalent standard SQL command.

See Also

[BEGIN](#), [COMMIT](#), [ROLLBACK](#)

ALTER AGGREGATE

ALTER AGGREGATE — change the definition of an aggregate function

Synopsis

```
ALTER AGGREGATE name ( aggregate_signature ) RENAME TO new_name
ALTER AGGREGATE name ( aggregate_signature )
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
ALTER AGGREGATE name ( aggregate_signature ) SET SCHEMA new_schema
```

where *aggregate_signature* is:

```
* |
[ argmode ] [ argname ] argtype [ , ... ] |
[ [ argmode ] [ argname ] argtype [ , ... ] ] ORDER BY [ argmode ] [ argname ] argtype
[ , ... ]
```

Description

ALTER AGGREGATE changes the definition of an aggregate function.

You must own the aggregate function to use ALTER AGGREGATE. To change the schema of an aggregate function, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the aggregate function's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the aggregate function. However, a superuser can alter ownership of any aggregate function anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing aggregate function.

argmode

The mode of an argument: IN or VARIADIC. If omitted, the default is IN.

argname

The name of an argument. Note that ALTER AGGREGATE does not actually pay any attention to argument names, since only the argument data types are needed to determine the aggregate function's identity.

argtype

An input data type on which the aggregate function operates. To reference a zero-argument aggregate function, write * in place of the list of argument specifications. To reference an ordered-set aggregate function, write ORDER BY between the direct and aggregated argument specifications.

new_name

The new name of the aggregate function.

new_owner

The new owner of the aggregate function.

new_schema

The new schema for the aggregate function.

Notes

The recommended syntax for referencing an ordered-set aggregate is to write `ORDER BY` between the direct and aggregated argument specifications, in the same style as in [CREATE AGGREGATE](#). However, it will also work to omit `ORDER BY` and just run the direct and aggregated argument specifications into a single list. In this abbreviated form, if `VARIADIC "any"` was used in both the direct and aggregated argument lists, write `VARIADIC "any"` only once.

Examples

To rename the aggregate function `myavg` for type `integer` to `my_average`:

```
ALTER AGGREGATE myavg(integer) RENAME TO my_average;
```

To change the owner of the aggregate function `myavg` for type `integer` to `joe`:

```
ALTER AGGREGATE myavg(integer) OWNER TO joe;
```

To move the ordered-set aggregate `mypercentile` with direct argument of type `float8` and aggregated argument of type `integer` into schema `myschema`:

```
ALTER AGGREGATE mypercentile(float8 ORDER BY integer) SET SCHEMA myschema;
```

This will work too:

```
ALTER AGGREGATE mypercentile(float8, integer) SET SCHEMA myschema;
```

Compatibility

There is no `ALTER AGGREGATE` statement in the SQL standard.

See Also

[CREATE AGGREGATE](#), [DROP AGGREGATE](#)

ALTER COLLATION

ALTER COLLATION — change the definition of a collation

Synopsis

```
ALTER COLLATION name REFRESH VERSION

ALTER COLLATION name RENAME TO new_name
ALTER COLLATION name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
    SESSION_USER }
ALTER COLLATION name SET SCHEMA new_schema
```

Description

ALTER COLLATION changes the definition of a collation.

You must own the collation to use ALTER COLLATION. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the collation's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the collation. However, a superuser can alter ownership of any collation anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing collation.

new_name

The new name of the collation.

new_owner

The new owner of the collation.

new_schema

The new schema for the collation.

REFRESH VERSION

Update the collation's version. See [Notes](#) below.

Notes

When a collation object is created, the provider-specific version of the collation is recorded in the system catalog. When the collation is used, the current version is checked against the recorded version, and a warning is issued when there is a mismatch, for example:

```
WARNING: collation "xx-x-icu" has version mismatch
DETAIL: The collation in the database was created using version 1.2.3.4, but the
operating system provides version 2.3.4.5.
HINT: Rebuild all objects affected by this collation and run ALTER COLLATION
pg_catalog."xx-x-icu" REFRESH VERSION, or build Postgres Pro with the right library
version.
```

A change in collation definitions can lead to corrupt indexes and other problems because the database system relies on stored objects having a certain sort order. Generally, this should be avoided, but it can happen in legitimate circumstances, such as when upgrading the operating system to a new major version or when using pg_upgrade to upgrade to server binaries linked with a newer version of ICU.

When this happens, all objects depending on the collation should be rebuilt, for example, using `REINDEX`. When that is done, the collation version can be refreshed using the command `ALTER COLLATION ... REFRESH VERSION`. This will update the system catalog to record the current collation version and will make the warning go away. Note that this does not actually check whether all affected objects have been rebuilt correctly.

When using collations provided by `libc`, version information is recorded on systems using the GNU C library (most Linux systems), FreeBSD and Windows. When using collations provided by ICU, the version information is provided by the ICU library and is available on all platforms.

Note

When using the GNU C library for collations, the C library's version is used as a proxy for the collation version. Many Linux distributions change collation definitions only when upgrading the C library, but this approach is imperfect as maintainers are free to back-port newer collation definitions to older C library releases.

When using Windows for collations, version information is only available for collations defined with BCP 47 language tags such as `en-US`.

For the database default collation, there is an analogous command `ALTER DATABASE ... REFRESH COLLATION VERSION`.

The following query can be used to identify all collations in the current database that need to be refreshed and the objects that depend on them:

```
SELECT pg_describe_object(refclassid, refobjid, refobjsubid) AS "Collation",
       pg_describe_object(classid, objid, objsubid) AS "Object"
FROM   pg_depend d JOIN pg_collation c
       ON refclassid = 'pg_collation'::regclass AND refobjid = c.oid
WHERE  c.collversion <> pg_collation_actual_version(c.oid)
ORDER BY 1, 2;
```

Examples

To rename the collation `ru_RU` to `russian`:

```
ALTER COLLATION "ru_RU" RENAME TO russian;
```

To change the owner of the collation `en_US` to `joe`:

```
ALTER COLLATION "en_US" OWNER TO joe;
```

Compatibility

There is no `ALTER COLLATION` statement in the SQL standard.

See Also

[CREATE COLLATION](#), [DROP COLLATION](#)

ALTER CONVERSION

ALTER CONVERSION — change the definition of a conversion

Synopsis

```
ALTER CONVERSION name RENAME TO new_name
ALTER CONVERSION name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER CONVERSION name SET SCHEMA new_schema
```

Description

ALTER CONVERSION changes the definition of a conversion.

You must own the conversion to use ALTER CONVERSION. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the conversion's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the conversion. However, a superuser can alter ownership of any conversion anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing conversion.

new_name

The new name of the conversion.

new_owner

The new owner of the conversion.

new_schema

The new schema for the conversion.

Examples

To rename the conversion `iso_8859_1_to_utf8` to `latin1_to_unicode`:

```
ALTER CONVERSION iso_8859_1_to_utf8 RENAME TO latin1_to_unicode;
```

To change the owner of the conversion `iso_8859_1_to_utf8` to `joe`:

```
ALTER CONVERSION iso_8859_1_to_utf8 OWNER TO joe;
```

Compatibility

There is no ALTER CONVERSION statement in the SQL standard.

See Also

[CREATE CONVERSION](#), [DROP CONVERSION](#)

ALTER DATABASE

ALTER DATABASE — change a database

Synopsis

```
ALTER DATABASE name [ [ WITH ] option [ ... ] ]
```

where *option* can be:

```
ALLOW_CONNECTIONS allowconn  
CONNECTION LIMIT connlimit  
IS_TEMPLATE istemplate
```

```
ALTER DATABASE name RENAME TO new_name
```

```
ALTER DATABASE name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

```
ALTER DATABASE name SET TABLESPACE new_tablespace
```

```
ALTER DATABASE name REFRESH COLLATION VERSION
```

```
ALTER DATABASE name SET configuration_parameter { TO | = } { value | DEFAULT }
```

```
ALTER DATABASE name SET configuration_parameter FROM CURRENT
```

```
ALTER DATABASE name RESET configuration_parameter
```

```
ALTER DATABASE name RESET ALL
```

Description

ALTER DATABASE changes the attributes of a database.

The first form changes certain per-database settings. (See below for details.) Only the database owner or a superuser can change these settings.

The second form changes the name of the database. Only the database owner or a superuser can rename a database; non-superuser owners must also have the `CREATEDB` privilege. The current database cannot be renamed. (Connect to a different database if you need to do that.)

The third form changes the owner of the database. To alter the owner, you must be able to `SET ROLE` to the new owning role, and you must have the `CREATEDB` privilege. (Note that superusers have all these privileges automatically.)

The fourth form changes the default tablespace of the database. Only the database owner or a superuser can do this; you must also have create privilege for the new tablespace. This command physically moves any tables or indexes in the database's old default tablespace to the new tablespace. The new default tablespace must be empty for this database, and no one can be connected to the database. Tables and indexes in non-default tablespaces are unaffected.

The remaining forms change the session default for a run-time configuration variable for a PostgreSQL database. Whenever a new session is subsequently started in that database, the specified value becomes the session default value. The database-specific default overrides whatever setting is present in `postgresql.conf` or has been received from the `postgres` command line. Only the database owner or a superuser can change the session defaults for a database. Certain variables cannot be set this way, or can only be set by a superuser.

Parameters

name

The name of the database whose attributes are to be altered.

allowconn

If false then no one can connect to this database.

conndefault

How many concurrent connections can be made to this database. -1 means no limit.

istemplate

If true, then this database can be cloned by any user with `CREATEDB` privileges; if false, then only superusers or the owner of the database can clone it.

new_name

The new name of the database.

new_owner

The new owner of the database.

new_tablespace

The new default tablespace of the database.

This form of the command cannot be executed inside a transaction block.

`REFRESH COLLATION VERSION`

Update the database collation version. See [Notes](#) for background.

configuration_parameter

value

Set this database's session default for the specified configuration parameter to the given value. If *value* is `DEFAULT` or, equivalently, `RESET` is used, the database-specific setting is removed, so the system-wide default setting will be inherited in new sessions. Use `RESET ALL` to clear all database-specific settings. `SET FROM CURRENT` saves the session's current value of the parameter as the database-specific value.

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

Notes

It is also possible to tie a session default to a specific role rather than to a database; see [ALTER ROLE](#). Role-specific settings override database-specific ones if there is a conflict.

Examples

To disable index scans by default in the database `test`:

```
ALTER DATABASE test SET enable_indexscan TO off;
```

Compatibility

The `ALTER DATABASE` statement is a Postgres Pro extension.

See Also

[CREATE DATABASE](#), [DROP DATABASE](#), [SET](#), [CREATE TABLESPACE](#)

ALTER DEFAULT PRIVILEGES

ALTER DEFAULT PRIVILEGES — define default access privileges

Synopsis

```
ALTER DEFAULT PRIVILEGES
    [ FOR { ROLE | USER } target_role [, ...] ]
    [ IN SCHEMA schema_name [, ...] ]
    abbreviated_grant_or_revoke
```

where *abbreviated_grant_or_revoke* is one of:

```
GRANT { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER |
        MAINTAIN }
        [, ...] | ALL [ PRIVILEGES ] }
ON TABLES
TO { [ GROUP ] role_name | PUBLIC } [, ...] [ WITH GRANT OPTION ]
```

```
GRANT { { USAGE | SELECT | UPDATE }
        [, ...] | ALL [ PRIVILEGES ] }
ON SEQUENCES
TO { [ GROUP ] role_name | PUBLIC } [, ...] [ WITH GRANT OPTION ]
```

```
GRANT { EXECUTE | ALL [ PRIVILEGES ] }
ON { FUNCTIONS | ROUTINES }
TO { [ GROUP ] role_name | PUBLIC } [, ...] [ WITH GRANT OPTION ]
```

```
GRANT { USAGE | ALL [ PRIVILEGES ] }
ON TYPES
TO { [ GROUP ] role_name | PUBLIC } [, ...] [ WITH GRANT OPTION ]
```

```
GRANT { { USAGE | CREATE }
        [, ...] | ALL [ PRIVILEGES ] }
ON SCHEMAS
TO { [ GROUP ] role_name | PUBLIC } [, ...] [ WITH GRANT OPTION ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER |
      MAINTAIN }
      [, ...] | ALL [ PRIVILEGES ] }
ON TABLES
FROM { [ GROUP ] role_name | PUBLIC } [, ...]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { { USAGE | SELECT | UPDATE }
      [, ...] | ALL [ PRIVILEGES ] }
ON SEQUENCES
FROM { [ GROUP ] role_name | PUBLIC } [, ...]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { EXECUTE | ALL [ PRIVILEGES ] }
ON { FUNCTIONS | ROUTINES }
FROM { [ GROUP ] role_name | PUBLIC } [, ...]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { USAGE | ALL [ PRIVILEGES ] }
    ON TYPES
    FROM { [ GROUP ] role_name | PUBLIC } [, ...]
    [ CASCADE | RESTRICT ]

REVOKE [ GRANT OPTION FOR ]
    { { USAGE | CREATE }
    [, ...] | ALL [ PRIVILEGES ] }
    ON SCHEMAS
    FROM { [ GROUP ] role_name | PUBLIC } [, ...]
    [ CASCADE | RESTRICT ]
```

Description

`ALTER DEFAULT PRIVILEGES` allows you to set the privileges that will be applied to objects created in the future. (It does not affect privileges assigned to already-existing objects.) Privileges can be set globally (i.e., for all objects created in the current database), or just for objects created in specified schemas.

While you can change your own default privileges and the defaults of roles that you are a member of, at object creation time, new object permissions are only affected by the default privileges of the current role, and are not inherited from any roles in which the current role is a member.

As explained in [Section 5.8](#), the default privileges for any object type normally grant all grantable permissions to the object owner, and may grant some privileges to `PUBLIC` as well. However, this behavior can be changed by altering the global default privileges with `ALTER DEFAULT PRIVILEGES`.

Currently, only the privileges for schemas, tables (including views and foreign tables), sequences, functions, and types (including domains) can be altered. For this command, functions include aggregates and procedures. The words `FUNCTIONS` and `ROUTINES` are equivalent in this command. (`ROUTINES` is preferred going forward as the standard term for functions and procedures taken together. In earlier Postgres Pro releases, only the word `FUNCTIONS` was allowed. It is not possible to set default privileges for functions and procedures separately.)

Default privileges that are specified per-schema are added to whatever the global default privileges are for the particular object type. This means you cannot revoke privileges per-schema if they are granted globally (either by default, or according to a previous `ALTER DEFAULT PRIVILEGES` command that did not specify a schema). Per-schema `REVOKE` is only useful to reverse the effects of a previous per-schema `GRANT`.

Parameters

target_role

Change default privileges for objects created by the *target_role*, or the current role if unspecified.

schema_name

The name of an existing schema. If specified, the default privileges are altered for objects later created in that schema. If `IN SCHEMA` is omitted, the global default privileges are altered. `IN SCHEMA` is not allowed when setting privileges for schemas, since schemas can't be nested.

role_name

The name of an existing role to grant or revoke privileges for. This parameter, and all the other parameters in *abbreviated_grant_or_revoke*, act as described under [GRANT](#) or [REVOKE](#), except that one is setting permissions for a whole class of objects rather than specific named objects.

Notes

Use `psql`'s `\ddp` command to obtain information about existing assignments of default privileges. The meaning of the privilege display is the same as explained for `\dp` in [Section 5.8](#).

If you wish to drop a role for which the default privileges have been altered, it is necessary to reverse the changes in its default privileges or use `DROP OWNED BY` to get rid of the default privileges entry for the role.

Examples

Grant `SELECT` privilege to everyone for all tables (and views) you subsequently create in schema `myschema`, and allow role `webuser` to `INSERT` into them too:

```
ALTER DEFAULT PRIVILEGES IN SCHEMA myschema GRANT SELECT ON TABLES TO PUBLIC;  
ALTER DEFAULT PRIVILEGES IN SCHEMA myschema GRANT INSERT ON TABLES TO webuser;
```

Undo the above, so that subsequently-created tables won't have any more permissions than normal:

```
ALTER DEFAULT PRIVILEGES IN SCHEMA myschema REVOKE SELECT ON TABLES FROM PUBLIC;  
ALTER DEFAULT PRIVILEGES IN SCHEMA myschema REVOKE INSERT ON TABLES FROM webuser;
```

Remove the public `EXECUTE` permission that is normally granted on functions, for all functions subsequently created by role `admin`:

```
ALTER DEFAULT PRIVILEGES FOR ROLE admin REVOKE EXECUTE ON FUNCTIONS FROM PUBLIC;
```

Note however that you *cannot* accomplish that effect with a command limited to a single schema. This command has no effect, unless it is undoing a matching `GRANT`:

```
ALTER DEFAULT PRIVILEGES IN SCHEMA public REVOKE EXECUTE ON FUNCTIONS FROM PUBLIC;
```

That's because per-schema default privileges can only add privileges to the global setting, not remove privileges granted by it.

Compatibility

There is no `ALTER DEFAULT PRIVILEGES` statement in the SQL standard.

See Also

[GRANT](#), [REVOKE](#)

ALTER DOMAIN

ALTER DOMAIN — change the definition of a domain

Synopsis

```
ALTER DOMAIN name
    { SET DEFAULT expression | DROP DEFAULT }
ALTER DOMAIN name
    { SET | DROP } NOT NULL
ALTER DOMAIN name
    ADD domain_constraint [ NOT VALID ]
ALTER DOMAIN name
    DROP CONSTRAINT [ IF EXISTS ] constraint_name [ RESTRICT | CASCADE ]
ALTER DOMAIN name
    RENAME CONSTRAINT constraint_name TO new_constraint_name
ALTER DOMAIN name
    VALIDATE CONSTRAINT constraint_name
ALTER DOMAIN name
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
ALTER DOMAIN name
    RENAME TO new_name
ALTER DOMAIN name
    SET SCHEMA new_schema
```

where *domain_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ NOT NULL | CHECK (expression) }
```

Description

ALTER DOMAIN changes the definition of an existing domain. There are several sub-forms:

SET/DROP DEFAULT

These forms set or remove the default value for a domain. Note that defaults only apply to subsequent INSERT commands; they do not affect rows already in a table using the domain.

SET/DROP NOT NULL

These forms change whether a domain is marked to allow NULL values or to reject NULL values. You can only SET NOT NULL when the columns using the domain contain no null values.

ADD *domain_constraint* [NOT VALID]

This form adds a new constraint to a domain. When a new constraint is added to a domain, all columns using that domain will be checked against the newly added constraint. These checks can be suppressed by adding the new constraint using the NOT VALID option; the constraint can later be made valid using ALTER DOMAIN ... VALIDATE CONSTRAINT. Newly inserted or updated rows are always checked against all constraints, even those marked NOT VALID. NOT VALID is only accepted for CHECK constraints.

DROP CONSTRAINT [IF EXISTS]

This form drops constraints on a domain. If IF EXISTS is specified and the constraint does not exist, no error is thrown. In this case a notice is issued instead.

RENAME CONSTRAINT

This form changes the name of a constraint on a domain.

VALIDATE CONSTRAINT

This form validates a constraint previously added as `NOT VALID`, that is, it verifies that all values in table columns of the domain type satisfy the specified constraint.

OWNER

This form changes the owner of the domain to the specified user.

RENAME

This form changes the name of the domain.

SET SCHEMA

This form changes the schema of the domain. Any constraints associated with the domain are moved into the new schema as well.

You must own the domain to use `ALTER DOMAIN`. To change the schema of a domain, you must also have `CREATE` privilege on the new schema. To alter the owner, you must be able to `SET ROLE` to the new owning role, and that role must have `CREATE` privilege on the domain's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the domain. However, a superuser can alter ownership of any domain anyway.)

Parameters

name

The name (possibly schema-qualified) of an existing domain to alter.

domain_constraint

New domain constraint for the domain.

constraint_name

Name of an existing constraint to drop or rename.

NOT VALID

Do not verify existing stored data for constraint validity.

CASCADE

Automatically drop objects that depend on the constraint, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the constraint if there are any dependent objects. This is the default behavior.

new_name

The new name for the domain.

new_constraint_name

The new name for the constraint.

new_owner

The user name of the new owner of the domain.

new_schema

The new schema for the domain.

Notes

Although `ALTER DOMAIN ADD CONSTRAINT` attempts to verify that existing stored data satisfies the new constraint, this check is not bulletproof, because the command cannot “see” table rows that are newly inserted or updated and not yet committed. If there is a hazard that concurrent operations might insert bad data, the way to proceed is to add the constraint using the `NOT VALID` option, commit that command, wait until all transactions started before that commit have finished, and then issue `ALTER DOMAIN VALIDATE CONSTRAINT` to search for data violating the constraint. This method is reliable because once the constraint is committed, all new transactions are guaranteed to enforce it against new values of the domain type.

Currently, `ALTER DOMAIN ADD CONSTRAINT`, `ALTER DOMAIN VALIDATE CONSTRAINT`, and `ALTER DOMAIN SET NOT NULL` will fail if the named domain or any derived domain is used within a container-type column (a composite, array, or range column) in any table in the database. They should eventually be improved to be able to verify the new constraint for such nested values.

Examples

To add a `NOT NULL` constraint to a domain:

```
ALTER DOMAIN zipcode SET NOT NULL;
```

To remove a `NOT NULL` constraint from a domain:

```
ALTER DOMAIN zipcode DROP NOT NULL;
```

To add a check constraint to a domain:

```
ALTER DOMAIN zipcode ADD CONSTRAINT zipchk CHECK (char_length(VALUE) = 5);
```

To remove a check constraint from a domain:

```
ALTER DOMAIN zipcode DROP CONSTRAINT zipchk;
```

To rename a check constraint on a domain:

```
ALTER DOMAIN zipcode RENAME CONSTRAINT zipchk TO zip_check;
```

To move the domain into a different schema:

```
ALTER DOMAIN zipcode SET SCHEMA customers;
```

Compatibility

`ALTER DOMAIN` conforms to the SQL standard, except for the `OWNER`, `RENAME`, `SET SCHEMA`, and `VALIDATE CONSTRAINT` variants, which are Postgres Pro extensions. The `NOT VALID` clause of the `ADD CONSTRAINT` variant is also a Postgres Pro extension.

See Also

[CREATE DOMAIN](#), [DROP DOMAIN](#)

ALTER EVENT TRIGGER

ALTER EVENT TRIGGER — change the definition of an event trigger

Synopsis

```
ALTER EVENT TRIGGER name DISABLE
ALTER EVENT TRIGGER name ENABLE [ REPLICA | ALWAYS ]
ALTER EVENT TRIGGER name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER EVENT TRIGGER name RENAME TO new_name
```

Description

ALTER EVENT TRIGGER changes properties of an existing event trigger.

You must be superuser to alter an event trigger.

Parameters

name

The name of an existing trigger to alter.

new_owner

The user name of the new owner of the event trigger.

new_name

The new name of the event trigger.

DISABLE/ENABLE [REPLICA | ALWAYS]

These forms configure the firing of event triggers. A disabled trigger is still known to the system, but is not executed when its triggering event occurs. See also [session_replication_role](#).

Compatibility

There is no ALTER EVENT TRIGGER statement in the SQL standard.

See Also

[CREATE EVENT TRIGGER](#), [DROP EVENT TRIGGER](#)

ALTER EXTENSION

ALTER EXTENSION — change the definition of an extension

Synopsis

```
ALTER EXTENSION name UPDATE [ TO new_version ]
ALTER EXTENSION name SET SCHEMA new_schema
ALTER EXTENSION name ADD member_object
ALTER EXTENSION name DROP member_object
```

where *member_object* is:

```
ACCESS METHOD object_name |
AGGREGATE aggregate_name ( aggregate_signature ) |
CAST ( source_type AS target_type ) |
COLLATION object_name |
CONVERSION object_name |
DOMAIN object_name |
EVENT TRIGGER object_name |
FOREIGN DATA WRAPPER object_name |
FOREIGN TABLE object_name |
FUNCTION function_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
MATERIALIZED VIEW object_name |
OPERATOR operator_name ( left_type, right_type ) |
OPERATOR CLASS object_name USING index_method |
OPERATOR FAMILY object_name USING index_method |
[ PROCEDURAL ] LANGUAGE object_name |
PROCEDURE procedure_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
ROUTINE routine_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
SCHEMA object_name |
SEQUENCE object_name |
SERVER object_name |
TABLE object_name |
TEXT SEARCH CONFIGURATION object_name |
TEXT SEARCH DICTIONARY object_name |
TEXT SEARCH PARSER object_name |
TEXT SEARCH TEMPLATE object_name |
TRANSFORM FOR type_name LANGUAGE lang_name |
TYPE object_name |
VIEW object_name
```

and *aggregate_signature* is:

```
* |
[ argmode ] [ argname ] argtype [ , ... ] |
[ [ argmode ] [ argname ] argtype [ , ... ] ] ORDER BY [ argmode ] [ argname ] argtype
[ , ... ]
```

Description

ALTER EXTENSION changes the definition of an installed extension. There are several subforms:

UPDATE

This form updates the extension to a newer version. The extension must supply a suitable update script (or series of scripts) that can modify the currently-installed version into the requested version.

SET SCHEMA

This form moves the extension's objects into another schema. The extension has to be *relocatable* for this command to succeed.

ADD *member_object*

This form adds an existing object to the extension. This is mainly useful in extension update scripts. The object will subsequently be treated as a member of the extension; notably, it can only be dropped by dropping the extension.

DROP *member_object*

This form removes a member object from the extension. This is mainly useful in extension update scripts. The object is not dropped, only disassociated from the extension.

See [Section 39.17](#) for more information about these operations.

You must own the extension to use `ALTER EXTENSION`. The `ADD/DROP` forms require ownership of the added/dropped object as well.

Parameters

name

The name of an installed extension.

new_version

The desired new version of the extension. This can be written as either an identifier or a string literal. If not specified, `ALTER EXTENSION UPDATE` attempts to update to whatever is shown as the default version in the extension's control file.

new_schema

The new schema for the extension.

*object_name**aggregate_name**function_name**operator_name**procedure_name**routine_name*

The name of an object to be added to or removed from the extension. Names of tables, aggregates, domains, foreign tables, functions, operators, operator classes, operator families, procedures, routines, sequences, text search objects, types, and views can be schema-qualified.

source_type

The name of the source data type of the cast.

target_type

The name of the target data type of the cast.

argmode

The mode of a function, procedure, or aggregate argument: `IN`, `OUT`, `INOUT`, or `VARIADIC`. If omitted, the default is `IN`. Note that `ALTER EXTENSION` does not actually pay any attention to `OUT` arguments, since only the input arguments are needed to determine the function's identity. So it is sufficient to list the `IN`, `INOUT`, and `VARIADIC` arguments.

argname

The name of a function, procedure, or aggregate argument. Note that `ALTER EXTENSION` does not actually pay any attention to argument names, since only the argument data types are needed to determine the function's identity.

argtype

The data type of a function, procedure, or aggregate argument.

left_type

right_type

The data type(s) of the operator's arguments (optionally schema-qualified). Write `NONE` for the missing argument of a prefix operator.

`PROCEDURAL`

This is a noise word.

type_name

The name of the data type of the transform.

lang_name

The name of the language of the transform.

Examples

To update the `hstore` extension to version 2.0:

```
ALTER EXTENSION hstore UPDATE TO '2.0';
```

To change the schema of the `hstore` extension to `utils`:

```
ALTER EXTENSION hstore SET SCHEMA utils;
```

To add an existing function to the `hstore` extension:

```
ALTER EXTENSION hstore ADD FUNCTION populate_record(anelement, hstore);
```

Compatibility

`ALTER EXTENSION` is a Postgres Pro extension.

See Also

[CREATE EXTENSION](#), [DROP EXTENSION](#)

ALTER FUNCTION

ALTER FUNCTION — change the definition of a function

Synopsis

```
ALTER FUNCTION name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    action [ ... ] [ RESTRICT ]  
ALTER FUNCTION name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    RENAME TO new_name  
ALTER FUNCTION name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }  
ALTER FUNCTION name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    SET SCHEMA new_schema  
ALTER FUNCTION name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    [ NO ] DEPENDS ON EXTENSION extension_name
```

where *action* is one of:

```
    CALLED ON NULL INPUT | RETURNS NULL ON NULL INPUT | STRICT  
    IMMUTABLE | STABLE | VOLATILE  
    [ NOT ] LEAKPROOF  
    [ EXTERNAL ] SECURITY INVOKER | [ EXTERNAL ] SECURITY DEFINER  
    PARALLEL { UNSAFE | RESTRICTED | SAFE }  
    COST execution_cost  
    ROWS result_rows  
    SUPPORT support_function  
    SET configuration_parameter { TO | = } { value | DEFAULT }  
    SET configuration_parameter FROM CURRENT  
    RESET configuration_parameter  
    RESET ALL
```

Description

ALTER FUNCTION changes the definition of a function.

You must own the function to use ALTER FUNCTION. To change a function's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the function's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the function. However, a superuser can alter ownership of any function anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing function. If no argument list is specified, the name must be unique in its schema.

argmode

The mode of an argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN. Note that ALTER FUNCTION does not actually pay any attention to OUT arguments, since only the input arguments are needed to determine the function's identity. So it is sufficient to list the IN, INOUT, and VARIADIC arguments.

argname

The name of an argument. Note that ALTER FUNCTION does not actually pay any attention to argument names, since only the argument data types are needed to determine the function's identity.

argtype

The data type(s) of the function's arguments (optionally schema-qualified), if any.

new_name

The new name of the function.

new_owner

The new owner of the function. Note that if the function is marked `SECURITY DEFINER`, it will subsequently execute as the new owner.

new_schema

The new schema for the function.

`DEPENDS ON EXTENSION extension_name`

`NO DEPENDS ON EXTENSION extension_name`

This form marks the function as dependent on the extension, or no longer dependent on that extension if `NO` is specified. A function that's marked as dependent on an extension is dropped when the extension is dropped, even if `CASCADE` is not specified. A function can depend upon multiple extensions, and will be dropped when any one of those extensions is dropped.

`CALLED ON NULL INPUT`

`RETURNS NULL ON NULL INPUT`

`STRICT`

`CALLED ON NULL INPUT` changes the function so that it will be invoked when some or all of its arguments are null. `RETURNS NULL ON NULL INPUT` or `STRICT` changes the function so that it is not invoked if any of its arguments are null; instead, a null result is assumed automatically. See [CREATE FUNCTION](#) for more information.

`IMMUTABLE`

`STABLE`

`VOLATILE`

Change the volatility of the function to the specified setting. See [CREATE FUNCTION](#) for details.

`[EXTERNAL] SECURITY INVOKER`

`[EXTERNAL] SECURITY DEFINER`

Change whether the function is a security definer or not. The key word `EXTERNAL` is ignored for SQL conformance. See [CREATE FUNCTION](#) for more information about this capability.

`PARALLEL`

Change whether the function is deemed safe for parallelism. See [CREATE FUNCTION](#) for details.

`LEAKPROOF`

Change whether the function is considered leakproof or not. See [CREATE FUNCTION](#) for more information about this capability.

`COST execution_cost`

Change the estimated execution cost of the function. See [CREATE FUNCTION](#) for more information.

`ROWS result_rows`

Change the estimated number of rows returned by a set-returning function. See [CREATE FUNCTION](#) for more information.

`SUPPORT support_function`

Set or change the planner support function to use for this function. See [Section 39.11](#) for details. You must be superuser to use this option.

This option cannot be used to remove the support function altogether, since it must name a new support function. Use `CREATE OR REPLACE FUNCTION` if you need to do that.

configuration_parameter
value

Add or change the assignment to be made to a configuration parameter when the function is called. If *value* is `DEFAULT` or, equivalently, `RESET` is used, the function-local setting is removed, so that the function executes with the value present in its environment. Use `RESET ALL` to clear all function-local settings. `SET FROM CURRENT` saves the value of the parameter that is current when `ALTER FUNCTION` is executed as the value to be applied when the function is entered.

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

`RESTRICT`

Ignored for conformance with the SQL standard.

Examples

To rename the function `sqrt` for type `integer` to `square_root`:

```
ALTER FUNCTION sqrt(integer) RENAME TO square_root;
```

To change the owner of the function `sqrt` for type `integer` to `joe`:

```
ALTER FUNCTION sqrt(integer) OWNER TO joe;
```

To change the schema of the function `sqrt` for type `integer` to `maths`:

```
ALTER FUNCTION sqrt(integer) SET SCHEMA maths;
```

To mark the function `sqrt` for type `integer` as being dependent on the extension `mathlib`:

```
ALTER FUNCTION sqrt(integer) DEPENDS ON EXTENSION mathlib;
```

To adjust the search path that is automatically set for a function:

```
ALTER FUNCTION check_password(text) SET search_path = admin, pg_temp;
```

To disable automatic setting of `search_path` for a function:

```
ALTER FUNCTION check_password(text) RESET search_path;
```

The function will now execute with whatever search path is used by its caller.

Compatibility

This statement is partially compatible with the `ALTER FUNCTION` statement in the SQL standard. The standard allows more properties of a function to be modified, but does not provide the ability to rename a function, make a function a security definer, attach configuration parameter values to a function, or change the owner, schema, or volatility of a function. The standard also requires the `RESTRICT` key word, which is optional in Postgres Pro.

See Also

[CREATE FUNCTION](#), [DROP FUNCTION](#), [ALTER PROCEDURE](#), [ALTER ROUTINE](#)

ALTER GROUP

ALTER GROUP — change role name or membership

Synopsis

```
ALTER GROUP role_specification ADD USER user_name [, ... ]  
ALTER GROUP role_specification DROP USER user_name [, ... ]
```

where *role_specification* can be:

```
role_name  
| CURRENT_ROLE  
| CURRENT_USER  
| SESSION_USER
```

```
ALTER GROUP group_name RENAME TO new_name
```

Description

ALTER GROUP changes the attributes of a user group. This is an obsolete command, though still accepted for backwards compatibility, because groups (and users too) have been superseded by the more general concept of roles.

The first two variants add users to a group or remove them from a group. (Any role can play the part of either a “user” or a “group” for this purpose.) These variants are effectively equivalent to granting or revoking membership in the role named as the “group”; so the preferred way to do this is to use [GRANT](#) or [REVOKE](#). Note that GRANT and REVOKE have additional options which are not available with this command, such as the ability to grant and revoke ADMIN OPTION, and the ability to specify the grantor.

The third variant changes the name of the group. This is exactly equivalent to renaming the role with [ALTER ROLE](#).

Parameters

group_name

The name of the group (role) to modify.

user_name

Users (roles) that are to be added to or removed from the group. The users must already exist; ALTER GROUP does not create or drop users.

new_name

The new name of the group.

Examples

Add users to a group:

```
ALTER GROUP staff ADD USER karl, john;
```

Remove a user from a group:

```
ALTER GROUP workers DROP USER beth;
```

Compatibility

There is no ALTER GROUP statement in the SQL standard.

See Also

[GRANT](#), [REVOKE](#), [ALTER ROLE](#)

ALTER INDEX

ALTER INDEX — change the definition of an index

Synopsis

```
ALTER INDEX [ IF EXISTS ] name RENAME TO new_name
ALTER INDEX [ IF EXISTS ] name SET TABLESPACE tablespace_name
ALTER INDEX name ATTACH PARTITION index_name
ALTER INDEX name [ NO ] DEPENDS ON EXTENSION extension_name
ALTER INDEX [ IF EXISTS ] name SET ( storage_parameter [= value] [, ... ] )
ALTER INDEX [ IF EXISTS ] name RESET ( storage_parameter [, ... ] )
ALTER INDEX [ IF EXISTS ] name ALTER [ COLUMN ] column_number
    SET STATISTICS integer
ALTER INDEX ALL IN TABLESPACE name [ OWNED BY role_name [, ... ] ]
    SET TABLESPACE new_tablespace [ NOWAIT ]
```

Description

ALTER INDEX changes the definition of an existing index. There are several subforms described below. Note that the lock level required may differ for each subform. An ACCESS EXCLUSIVE lock is held unless explicitly noted. When multiple subcommands are listed, the lock held will be the strictest one required from any subcommand.

RENAME

The RENAME form changes the name of the index. If the index is associated with a table constraint (either UNIQUE, PRIMARY KEY, or EXCLUDE), the constraint is renamed as well. There is no effect on the stored data.

Renaming an index acquires a SHARE UPDATE EXCLUSIVE lock.

SET TABLESPACE

This form changes the index's tablespace to the specified tablespace and moves the data file(s) associated with the index to the new tablespace. To change the tablespace of an index, you must own the index and have CREATE privilege on the new tablespace. All indexes in the current database in a tablespace can be moved by using the ALL IN TABLESPACE form, which will lock all indexes to be moved and then move each one. This form also supports OWNED BY, which will only move indexes owned by the roles specified. If the NOWAIT option is specified then the command will fail if it is unable to acquire all of the locks required immediately. Note that system catalogs will not be moved by this command, use ALTER DATABASE or explicit ALTER INDEX invocations instead if desired. See also [CREATE TABLESPACE](#).

ATTACH PARTITION *index_name*

Causes the named index (possibly schema-qualified) to become attached to the altered index. The named index must be on a partition of the table containing the index being altered, and have an equivalent definition. An attached index cannot be dropped by itself, and will automatically be dropped if its parent index is dropped.

DEPENDS ON EXTENSION *extension_name*

NO DEPENDS ON EXTENSION *extension_name*

This form marks the index as dependent on the extension, or no longer dependent on that extension if NO is specified. An index that's marked as dependent on an extension is automatically dropped when the extension is dropped.

SET (*storage_parameter* [= *value*] [, ...])

This form changes one or more index-method-specific storage parameters for the index. See [CREATE INDEX](#) for details on the available parameters. Note that the index contents will not be modified

immediately by this command; depending on the parameter you might need to rebuild the index with [REINDEX](#) to get the desired effects.

```
RESET ( storage_parameter [, ... ] )
```

This form resets one or more index-method-specific storage parameters to their defaults. As with [SET](#), a [REINDEX](#) might be needed to update the index entirely.

```
ALTER [ COLUMN ] column_number SET STATISTICS integer
```

This form sets the per-column statistics-gathering target for subsequent [ANALYZE](#) operations, though can be used only on index columns that are defined as an expression. Since expressions lack a unique name, we refer to them using the ordinal number of the index column. The target can be set in the range 0 to 10000; alternatively, set it to -1 to revert to using the system default statistics target ([default_statistics_target](#)). For more information on the use of statistics by the Postgres Pro query planner, refer to [Section 14.2](#).

Parameters

IF EXISTS

Do not throw an error if the index does not exist. A notice is issued in this case.

column_number

The ordinal number refers to the ordinal (left-to-right) position of the index column.

name

The name (possibly schema-qualified) of an existing index to alter.

new_name

The new name for the index.

tablespace_name

The tablespace to which the index will be moved.

extension_name

The name of the extension that the index is to depend on.

storage_parameter

The name of an index-method-specific storage parameter.

value

The new value for an index-method-specific storage parameter. This might be a number or a word depending on the parameter.

Notes

These operations are also possible using [ALTER TABLE](#). `ALTER INDEX` is in fact just an alias for the forms of `ALTER TABLE` that apply to indexes.

There was formerly an `ALTER INDEX OWNER` variant, but this is now ignored (with a warning). An index cannot have an owner different from its table's owner. Changing the table's owner automatically changes the index as well.

Changing any part of a system catalog index is not permitted.

Examples

To rename an existing index:

```
ALTER INDEX distributors RENAME TO suppliers;
```

To move an index to a different tablespace:

```
ALTER INDEX distributors SET TABLESPACE fasttablespace;
```

To change an index's fill factor (assuming that the index method supports it):

```
ALTER INDEX distributors SET (fillfactor = 75);  
REINDEX INDEX distributors;
```

Set the statistics-gathering target for an expression index:

```
CREATE INDEX coord_idx ON measured (x, y, (z + t));  
ALTER INDEX coord_idx ALTER COLUMN 3 SET STATISTICS 1000;
```

Compatibility

ALTER INDEX is a Postgres Pro extension.

See Also

[CREATE INDEX](#), [REINDEX](#)

ALTER LANGUAGE

ALTER LANGUAGE — change the definition of a procedural language

Synopsis

```
ALTER [ PROCEDURAL ] LANGUAGE name RENAME TO new_name
ALTER [ PROCEDURAL ] LANGUAGE name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
```

Description

ALTER LANGUAGE changes the definition of a procedural language. The only functionality is to rename the language or assign a new owner. You must be superuser or owner of the language to use ALTER LANGUAGE.

Parameters

name

Name of a language

new_name

The new name of the language

new_owner

The new owner of the language

Compatibility

There is no ALTER LANGUAGE statement in the SQL standard.

See Also

[CREATE LANGUAGE](#), [DROP LANGUAGE](#)

ALTER LARGE OBJECT

ALTER LARGE OBJECT — change the definition of a large object

Synopsis

```
ALTER LARGE OBJECT large_object_oid OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER  
| SESSION_USER }
```

Description

ALTER LARGE OBJECT changes the definition of a large object.

You must own the large object to use ALTER LARGE OBJECT. To alter the owner, you must also be able to SET ROLE to the new owning role. (However, a superuser can alter any large object anyway.) Currently, the only functionality is to assign a new owner, so both restrictions always apply.

Parameters

large_object_oid

OID of the large object to be altered

new_owner

The new owner of the large object

Compatibility

There is no ALTER LARGE OBJECT statement in the SQL standard.

See Also

[Chapter 36](#)

ALTER MATERIALIZED VIEW

ALTER MATERIALIZED VIEW — change the definition of a materialized view

Synopsis

```
ALTER MATERIALIZED VIEW [ IF EXISTS ] name
    action [, ... ]
ALTER MATERIALIZED VIEW name
    [ NO ] DEPENDS ON EXTENSION extension_name
ALTER MATERIALIZED VIEW [ IF EXISTS ] name
    RENAME [ COLUMN ] column_name TO new_column_name
ALTER MATERIALIZED VIEW [ IF EXISTS ] name
    RENAME TO new_name
ALTER MATERIALIZED VIEW [ IF EXISTS ] name
    SET SCHEMA new_schema
ALTER MATERIALIZED VIEW ALL IN TABLESPACE name [ OWNED BY role_name [, ... ] ]
    SET TABLESPACE new_tablespace [ NOWAIT ]
```

where *action* is one of:

```
ALTER [ COLUMN ] column_name SET STATISTICS integer
ALTER [ COLUMN ] column_name SET ( attribute_option = value [, ... ] )
ALTER [ COLUMN ] column_name RESET ( attribute_option [, ... ] )
ALTER [ COLUMN ] column_name SET STORAGE { PLAIN | EXTERNAL | EXTENDED | MAIN |
DEFAULT }
ALTER [ COLUMN ] column_name SET COMPRESSION compression_method
CLUSTER ON index_name
SET WITHOUT CLUSTER
SET ACCESS METHOD new_access_method
SET TABLESPACE new_tablespace
SET ( storage_parameter [= value] [, ... ] )
RESET ( storage_parameter [, ... ] )
OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

Description

ALTER MATERIALIZED VIEW changes various auxiliary properties of an existing materialized view.

You must own the materialized view to use ALTER MATERIALIZED VIEW. To change a materialized view's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the materialized view's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the materialized view. However, a superuser can alter ownership of any view anyway.)

The statement subforms and actions available for ALTER MATERIALIZED VIEW are a subset of those available for ALTER TABLE, and have the same meaning when used for materialized views. See the descriptions for [ALTER TABLE](#) for details.

Parameters

name

The name (optionally schema-qualified) of an existing materialized view.

column_name

Name of an existing column.

extension_name

The name of the extension that the materialized view is to depend on (or no longer dependent on, if `NO` is specified). A materialized view that's marked as dependent on an extension is automatically dropped when the extension is dropped.

new_column_name

New name for an existing column.

new_owner

The user name of the new owner of the materialized view.

new_name

The new name for the materialized view.

new_schema

The new schema for the materialized view.

Examples

To rename the materialized view `foo` to `bar`:

```
ALTER MATERIALIZED VIEW foo RENAME TO bar;
```

Compatibility

`ALTER MATERIALIZED VIEW` is a Postgres Pro extension.

See Also

[CREATE MATERIALIZED VIEW](#), [DROP MATERIALIZED VIEW](#), [REFRESH MATERIALIZED VIEW](#)

ALTER OPERATOR

ALTER OPERATOR — change the definition of an operator

Synopsis

```
ALTER OPERATOR name ( { left_type | NONE } , right_type )  
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

```
ALTER OPERATOR name ( { left_type | NONE } , right_type )  
    SET SCHEMA new_schema
```

```
ALTER OPERATOR name ( { left_type | NONE } , right_type )  
    SET ( { RESTRICT = { res_proc | NONE }  
        | JOIN = { join_proc | NONE }  
        | COMMUTATOR = com_op  
        | NEGATOR = neg_op  
        | HASHES  
        | MERGES  
        } [, ... ] )
```

Description

ALTER OPERATOR changes the definition of an operator.

You must own the operator to use ALTER OPERATOR. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the operator's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the operator. However, a superuser can alter ownership of any operator anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing operator.

left_type

The data type of the operator's left operand; write NONE if the operator has no left operand.

right_type

The data type of the operator's right operand.

new_owner

The new owner of the operator.

new_schema

The new schema for the operator.

res_proc

The restriction selectivity estimator function for this operator; write NONE to remove existing selectivity estimator.

join_proc

The join selectivity estimator function for this operator; write NONE to remove existing selectivity estimator.

com_op

The commutator of this operator. Can only be changed if the operator does not have an existing commutator.

neg_op

The negator of this operator. Can only be changed if the operator does not have an existing negator.

HASHES

Indicates this operator can support a hash join. Can only be enabled and not disabled.

MERGES

Indicates this operator can support a merge join. Can only be enabled and not disabled.

Notes

Refer to [Section 39.14](#) and [Section 39.15](#) for further information.

Since commutators come in pairs that are commutators of each other, `ALTER OPERATOR SET COMMUTATOR` will also set the commutator of the *com_op* to be the target operator. Likewise, `ALTER OPERATOR SET NEGATOR` will also set the negator of the *neg_op* to be the target operator. Therefore, you must own the commutator or negator operator as well as the target operator.

Examples

Change the owner of a custom operator `a @@ b` for type `text`:

```
ALTER OPERATOR @@ (text, text) OWNER TO joe;
```

Change the restriction and join selectivity estimator functions of a custom operator `a && b` for type `int[]`:

```
ALTER OPERATOR && (int[], int[]) SET (RESTRICT = _int_contsel, JOIN =  
_int_contjoinssel);
```

Mark the `&&` operator as being its own commutator:

```
ALTER OPERATOR && (int[], int[]) SET (COMMUTATOR = &&);
```

Compatibility

There is no `ALTER OPERATOR` statement in the SQL standard.

See Also

[CREATE OPERATOR](#), [DROP OPERATOR](#)

ALTER OPERATOR CLASS

ALTER OPERATOR CLASS — change the definition of an operator class

Synopsis

```
ALTER OPERATOR CLASS name USING index_method
    RENAME TO new_name
```

```
ALTER OPERATOR CLASS name USING index_method
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

```
ALTER OPERATOR CLASS name USING index_method
    SET SCHEMA new_schema
```

Description

ALTER OPERATOR CLASS changes the definition of an operator class.

You must own the operator class to use ALTER OPERATOR CLASS. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the operator class's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the operator class. However, a superuser can alter ownership of any operator class anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing operator class.

index_method

The name of the index method this operator class is for.

new_name

The new name of the operator class.

new_owner

The new owner of the operator class.

new_schema

The new schema for the operator class.

Compatibility

There is no ALTER OPERATOR CLASS statement in the SQL standard.

See Also

[CREATE OPERATOR CLASS](#), [DROP OPERATOR CLASS](#), [ALTER OPERATOR FAMILY](#)

ALTER OPERATOR FAMILY

ALTER OPERATOR FAMILY — change the definition of an operator family

Synopsis

```
ALTER OPERATOR FAMILY name USING index_method ADD
{ OPERATOR strategy_number operator_name ( op_type, op_type )
  [ FOR SEARCH | FOR ORDER BY sort_family_name ]
| FUNCTION support_number [ ( op_type [ , op_type ] ) ]
  function_name [ ( argument_type [, ...] ) ]
} [, ... ]

ALTER OPERATOR FAMILY name USING index_method DROP
{ OPERATOR strategy_number ( op_type [ , op_type ] )
| FUNCTION support_number ( op_type [ , op_type ] )
} [, ... ]

ALTER OPERATOR FAMILY name USING index_method
  RENAME TO new_name

ALTER OPERATOR FAMILY name USING index_method
  OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }

ALTER OPERATOR FAMILY name USING index_method
  SET SCHEMA new_schema
```

Description

`ALTER OPERATOR FAMILY` changes the definition of an operator family. You can add operators and support functions to the family, remove them from the family, or change the family's name or owner.

When operators and support functions are added to a family with `ALTER OPERATOR FAMILY`, they are not part of any specific operator class within the family, but are just “loose” within the family. This indicates that these operators and functions are compatible with the family's semantics, but are not required for correct functioning of any specific index. (Operators and functions that are so required should be declared as part of an operator class, instead; see [CREATE OPERATOR CLASS](#).) Postgres Pro will allow loose members of a family to be dropped from the family at any time, but members of an operator class cannot be dropped without dropping the whole class and any indexes that depend on it. Typically, single-data-type operators and functions are part of operator classes because they are needed to support an index on that specific data type, while cross-data-type operators and functions are made loose members of the family.

You must be a superuser to use `ALTER OPERATOR FAMILY`. (This restriction is made because an erroneous operator family definition could confuse or even crash the server.)

`ALTER OPERATOR FAMILY` does not presently check whether the operator family definition includes all the operators and functions required by the index method, nor whether the operators and functions form a self-consistent set. It is the user's responsibility to define a valid operator family.

Refer to [Section 39.16](#) for further information.

Parameters

name

The name (optionally schema-qualified) of an existing operator family.

index_method

The name of the index method this operator family is for.

strategy_number

The index method's strategy number for an operator associated with the operator family.

operator_name

The name (optionally schema-qualified) of an operator associated with the operator family.

op_type

In an `OPERATOR` clause, the operand data type(s) of the operator, or `NONE` to signify a prefix operator. Unlike the comparable syntax in `CREATE OPERATOR CLASS`, the operand data types must always be specified.

In an `ADD FUNCTION` clause, the operand data type(s) the function is intended to support, if different from the input data type(s) of the function. For B-tree comparison functions and hash functions it is not necessary to specify *op_type* since the function's input data type(s) are always the correct ones to use. For B-tree sort support functions, B-Tree equal image functions, and all functions in GiST, SP-GiST and GIN operator classes, it is necessary to specify the operand data type(s) the function is to be used with.

In a `DROP FUNCTION` clause, the operand data type(s) the function is intended to support must be specified.

sort_family_name

The name (optionally schema-qualified) of an existing `btree` operator family that describes the sort ordering associated with an ordering operator.

If neither `FOR SEARCH` nor `FOR ORDER BY` is specified, `FOR SEARCH` is the default.

support_number

The index method's support function number for a function associated with the operator family.

function_name

The name (optionally schema-qualified) of a function that is an index method support function for the operator family. If no argument list is specified, the name must be unique in its schema.

argument_type

The parameter data type(s) of the function.

new_name

The new name of the operator family.

new_owner

The new owner of the operator family.

new_schema

The new schema for the operator family.

The `OPERATOR` and `FUNCTION` clauses can appear in any order.

Notes

Notice that the `DROP` syntax only specifies the “slot” in the operator family, by strategy or support number and input data type(s). The name of the operator or function occupying the slot is not mentioned. Also, for `DROP FUNCTION` the type(s) to specify are the input data type(s) the function is intended to support;

for GiST, SP-GiST and GIN indexes this might have nothing to do with the actual input argument types of the function.

Because the index machinery does not check access permissions on functions before using them, including a function or operator in an operator family is tantamount to granting public execute permission on it. This is usually not an issue for the sorts of functions that are useful in an operator family.

The operators should not be defined by SQL functions. An SQL function is likely to be inlined into the calling query, which will prevent the optimizer from recognizing that the query matches an index.

Before PostgreSQL 8.4, the `OPERATOR` clause could include a `RECHECK` option. This is no longer supported because whether an index operator is “lossy” is now determined on-the-fly at run time. This allows efficient handling of cases where an operator might or might not be lossy.

Examples

The following example command adds cross-data-type operators and support functions to an operator family that already contains B-tree operator classes for data types `int4` and `int2`.

```
ALTER OPERATOR FAMILY integer_ops USING btree ADD
```

```
-- int4 vs int2
OPERATOR 1 < (int4, int2) ,
OPERATOR 2 <= (int4, int2) ,
OPERATOR 3 = (int4, int2) ,
OPERATOR 4 >= (int4, int2) ,
OPERATOR 5 > (int4, int2) ,
FUNCTION 1 btint42cmp(int4, int2) ,

-- int2 vs int4
OPERATOR 1 < (int2, int4) ,
OPERATOR 2 <= (int2, int4) ,
OPERATOR 3 = (int2, int4) ,
OPERATOR 4 >= (int2, int4) ,
OPERATOR 5 > (int2, int4) ,
FUNCTION 1 btint24cmp(int2, int4) ;
```

To remove these entries again:

```
ALTER OPERATOR FAMILY integer_ops USING btree DROP
```

```
-- int4 vs int2
OPERATOR 1 (int4, int2) ,
OPERATOR 2 (int4, int2) ,
OPERATOR 3 (int4, int2) ,
OPERATOR 4 (int4, int2) ,
OPERATOR 5 (int4, int2) ,
FUNCTION 1 (int4, int2) ,

-- int2 vs int4
OPERATOR 1 (int2, int4) ,
OPERATOR 2 (int2, int4) ,
OPERATOR 3 (int2, int4) ,
OPERATOR 4 (int2, int4) ,
OPERATOR 5 (int2, int4) ,
FUNCTION 1 (int2, int4) ;
```

Compatibility

There is no `ALTER OPERATOR FAMILY` statement in the SQL standard.

See Also

[CREATE OPERATOR FAMILY](#), [DROP OPERATOR FAMILY](#), [CREATE OPERATOR CLASS](#), [ALTER OPERATOR CLASS](#), [DROP OPERATOR CLASS](#)

ALTER POLICY

ALTER POLICY — change the definition of a row-level security policy

Synopsis

```
ALTER POLICY name ON table_name RENAME TO new_name
```

```
ALTER POLICY name ON table_name  
[ TO { role_name | PUBLIC | CURRENT_ROLE | CURRENT_USER | SESSION_USER } [, ...] ]  
[ USING ( using_expression ) ]  
[ WITH CHECK ( check_expression ) ]
```

Description

ALTER POLICY changes the definition of an existing row-level security policy. Note that ALTER POLICY only allows the set of roles to which the policy applies and the USING and WITH CHECK expressions to be modified. To change other properties of a policy, such as the command to which it applies or whether it is permissive or restrictive, the policy must be dropped and recreated.

To use ALTER POLICY, you must own the table that the policy applies to.

In the second form of ALTER POLICY, the role list, *using_expression*, and *check_expression* are replaced independently if specified. When one of those clauses is omitted, the corresponding part of the policy is unchanged.

Parameters

name

The name of an existing policy to alter.

table_name

The name (optionally schema-qualified) of the table that the policy is on.

new_name

The new name for the policy.

role_name

The role(s) to which the policy applies. Multiple roles can be specified at one time. To apply the policy to all roles, use PUBLIC.

using_expression

The USING expression for the policy. See [CREATE POLICY](#) for details.

check_expression

The WITH CHECK expression for the policy. See [CREATE POLICY](#) for details.

Compatibility

ALTER POLICY is a Postgres Pro extension.

See Also

[CREATE POLICY](#), [DROP POLICY](#)

ALTER PROCEDURE

ALTER PROCEDURE — change the definition of a procedure

Synopsis

```
ALTER PROCEDURE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    action [ ... ] [ RESTRICT ]  
ALTER PROCEDURE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    RENAME TO new_name  
ALTER PROCEDURE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }  
ALTER PROCEDURE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    SET SCHEMA new_schema  
ALTER PROCEDURE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
    [ NO ] DEPENDS ON EXTENSION extension_name
```

where *action* is one of:

```
[ EXTERNAL ] SECURITY INVOKER | [ EXTERNAL ] SECURITY DEFINER  
SET configuration_parameter { TO | = } { value | DEFAULT }  
SET configuration_parameter FROM CURRENT  
RESET configuration_parameter  
RESET ALL
```

Description

ALTER PROCEDURE changes the definition of a procedure.

You must own the procedure to use ALTER PROCEDURE. To change a procedure's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the procedure's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the procedure. However, a superuser can alter ownership of any procedure anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing procedure. If no argument list is specified, the name must be unique in its schema.

argmode

The mode of an argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN.

argname

The name of an argument. Note that ALTER PROCEDURE does not actually pay any attention to argument names, since only the argument data types are used to determine the procedure's identity.

argtype

The data type(s) of the procedure's arguments (optionally schema-qualified), if any. See [DROP PROCEDURE](#) for the details of how the procedure is looked up using the argument data type(s).

new_name

The new name of the procedure.

new_owner

The new owner of the procedure. Note that if the procedure is marked `SECURITY DEFINER`, it will subsequently execute as the new owner.

new_schema

The new schema for the procedure.

extension_name

This form marks the procedure as dependent on the extension, or no longer dependent on the extension if `NO` is specified. A procedure that's marked as dependent on an extension is dropped when the extension is dropped, even if `cascade` is not specified. A procedure can depend upon multiple extensions, and will be dropped when any one of those extensions is dropped.

```
[ EXTERNAL ] SECURITY INVOKER  
[ EXTERNAL ] SECURITY DEFINER
```

Change whether the procedure is a security definer or not. The key word `EXTERNAL` is ignored for SQL conformance. See [CREATE PROCEDURE](#) for more information about this capability.

configuration_parameter
value

Add or change the assignment to be made to a configuration parameter when the procedure is called. If *value* is `DEFAULT` or, equivalently, `RESET` is used, the procedure-local setting is removed, so that the procedure executes with the value present in its environment. Use `RESET ALL` to clear all procedure-local settings. `SET FROM CURRENT` saves the value of the parameter that is current when `ALTER PROCEDURE` is executed as the value to be applied when the procedure is entered.

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

`RESTRICT`

Ignored for conformance with the SQL standard.

Examples

To rename the procedure `insert_data` with two arguments of type `integer` to `insert_record`:

```
ALTER PROCEDURE insert_data(integer, integer) RENAME TO insert_record;
```

To change the owner of the procedure `insert_data` with two arguments of type `integer` to `joe`:

```
ALTER PROCEDURE insert_data(integer, integer) OWNER TO joe;
```

To change the schema of the procedure `insert_data` with two arguments of type `integer` to `accounting`:

```
ALTER PROCEDURE insert_data(integer, integer) SET SCHEMA accounting;
```

To mark the procedure `insert_data(integer, integer)` as being dependent on the extension `myext`:

```
ALTER PROCEDURE insert_data(integer, integer) DEPENDS ON EXTENSION myext;
```

To adjust the search path that is automatically set for a procedure:

```
ALTER PROCEDURE check_password(text) SET search_path = admin, pg_temp;
```

To disable automatic setting of `search_path` for a procedure:

```
ALTER PROCEDURE check_password(text) RESET search_path;
```

The procedure will now execute with whatever search path is used by its caller.

Compatibility

This statement is partially compatible with the `ALTER PROCEDURE` statement in the SQL standard. The standard allows more properties of a procedure to be modified, but does not provide the ability to rename a procedure, make a procedure a security definer, attach configuration parameter values to a procedure, or change the owner, schema, or volatility of a procedure. The standard also requires the `RESTRICT` key word, which is optional in Postgres Pro.

See Also

[CREATE PROCEDURE](#), [DROP PROCEDURE](#), [ALTER FUNCTION](#), [ALTER ROUTINE](#)

ALTER PROFILE

ALTER PROFILE — change an authentication profile

Synopsis

```
ALTER PROFILE name [ LIMIT parameter value [ ... ] ]
```

where *parameter* can be:

```
    FAILED_LOGIN_ATTEMPTS
|  PASSWORD_REUSE_TIME
|  PASSWORD_REUSE_MAX
|  PASSWORD_LIFE_TIME
|  PASSWORD_GRACE_TIME
|  USER_INACTIVE_TIME
|  FAILED_AUTH_KEEP_TIME
|  PASSWORD_MIN_UNIQUE_CHARS
|  PASSWORD_MIN_LEN
|  PASSWORD_REQUIRE_COMPLEX
```

```
ALTER PROFILE name RENAME TO new_name
```

Description

The `ALTER PROFILE` command changes parameters of a Postgres Pro profile. You must be a database superuser or have the privileges of the [pg_manage_profiles](#) role to use this command.

Parameters

name

The name of the profile for which to alter parameters.

```
FAILED_LOGIN_ATTEMPTS value
PASSWORD_REUSE_TIME value
PASSWORD_REUSE_MAX value
PASSWORD_LIFE_TIME value
PASSWORD_GRACE_TIME value
USER_INACTIVE_TIME value
FAILED_AUTH_KEEP_TIME value
PASSWORD_MIN_UNIQUE_CHARS value
PASSWORD_MIN_LEN value
PASSWORD_REQUIRE_COMPLEX value
```

These clauses alter parameters originally set by [CREATE PROFILE](#). For more information, see the [CREATE PROFILE](#) reference page.

new_name

The new name of the profile.

Notes

Use [CREATE PROFILE](#) to add new profiles, and [DROP PROFILE](#) to remove a profile.

Examples

Change `PASSWORD_LIFE_TIME` and `PASSWORD_GRACE_TIME` parameters of a profile:

```
ALTER PROFILE admin_profile  
  LIMIT PASSWORD_LIFE_TIME 90  
  PASSWORD_GRACE_TIME 3;
```

See Also

[CREATE PROFILE](#), [DROP PROFILE](#), [CREATE ROLE](#)

ALTER PUBLICATION

ALTER PUBLICATION — change the definition of a publication

Synopsis

```
ALTER PUBLICATION name ADD publication_object [, ...]
ALTER PUBLICATION name SET publication_object [, ...]
ALTER PUBLICATION name DROP publication_object [, ...]
ALTER PUBLICATION name SET ( publication_parameter [= value] [, ... ] )
ALTER PUBLICATION name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
    SESSION_USER }
ALTER PUBLICATION name RENAME TO new_name
```

where *publication_object* is one of:

```
TABLE [ ONLY ] table_name [ * ] [ ( column_name [, ... ] ) ] [ WHERE ( expression
) ] [, ... ]
TABLES IN SCHEMA { schema_name | CURRENT_SCHEMA } [, ... ]
```

Description

The command ALTER PUBLICATION can change the attributes of a publication.

The first three variants change which tables/schemas are part of the publication. The SET clause will replace the list of tables/schemas in the publication with the specified list; the existing tables/schemas that were present in the publication will be removed. The ADD and DROP clauses will add and remove one or more tables/schemas from the publication. Note that adding tables/schemas to a publication that is already subscribed to will require an [ALTER SUBSCRIPTION ... REFRESH PUBLICATION](#) action on the subscribing side in order to become effective. Note also that DROP TABLES IN SCHEMA will not drop any schema tables that were specified using [FOR TABLE/ ADD TABLE](#), and the combination of DROP with a WHERE clause is not allowed.

The fourth variant of this command listed in the synopsis can change all of the publication properties specified in [CREATE PUBLICATION](#). Properties not mentioned in the command retain their previous settings.

The remaining variants change the owner and the name of the publication.

You must own the publication to use ALTER PUBLICATION. Adding a table to a publication additionally requires owning that table. The ADD TABLES IN SCHEMA and SET TABLES IN SCHEMA to a publication requires the invoking user to be a superuser. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the database. Also, the new owner of a [FOR ALL TABLES](#) or [FOR TABLES IN SCHEMA](#) publication must be a superuser. However, a superuser can change the ownership of a publication regardless of these restrictions.

Adding/Setting any schema when the publication also publishes a table with a column list, and vice versa is not supported.

Parameters

name

The name of an existing publication whose definition is to be altered.

table_name

Name of an existing table. If ONLY is specified before the table name, only that table is affected. If ONLY is not specified, the table and all its descendant tables (if any) are affected. Optionally, * can be specified after the table name to explicitly indicate that descendant tables are included.

Optionally, a column list can be specified. See [CREATE PUBLICATION](#) for details. Note that a subscription having several publications in which the same table has been published with different column lists is not supported. See [Warning: Combining Column Lists from Multiple Publications](#) for details of potential problems when altering column lists.

If the optional `WHERE` clause is specified, rows for which the *expression* evaluates to false or null will not be published. Note that parentheses are required around the expression. The *expression* is evaluated with the role used for the replication connection.

schema_name

Name of an existing schema.

`SET (publication_parameter [= value] [, ...])`

This clause alters publication parameters originally set by [CREATE PUBLICATION](#). See there for more information.

new_owner

The user name of the new owner of the publication.

new_name

The new name for the publication.

Examples

Change the publication to publish only deletes and updates:

```
ALTER PUBLICATION noinsert SET (publish = 'update, delete');
```

Add some tables to the publication:

```
ALTER PUBLICATION mypublication ADD TABLE users (user_id, firstname), departments;
```

Change the set of columns published for a table:

```
ALTER PUBLICATION mypublication SET TABLE users (user_id, firstname, lastname), TABLE departments;
```

Add schemas `marketing` and `sales` to the publication `sales_publication`:

```
ALTER PUBLICATION sales_publication ADD TABLES IN SCHEMA marketing, sales;
```

Add tables `users`, `departments` and schema `production` to the publication `production_publication`:

```
ALTER PUBLICATION production_publication ADD TABLE users, departments, TABLES IN SCHEMA production;
```

Compatibility

`ALTER PUBLICATION` is a Postgres Pro extension.

See Also

[CREATE PUBLICATION](#), [DROP PUBLICATION](#), [CREATE SUBSCRIPTION](#), [ALTER SUBSCRIPTION](#)

ALTER ROLE

ALTER ROLE — change a database role

Synopsis

```
ALTER ROLE role_specification [ WITH ] option [ ... ]
```

where *option* can be:

```
SUPERUSER | NOSUPERUSER
| CREATEDB | NOCREATEDB
| CREATEROLE | NOCREATEROLE
| INHERIT | NOINHERIT
| LOGIN | NOLOGIN
| REPLICATION | NOREPLICATION
| BYPASSRLS | NOBYPASSRLS
| CONNECTION LIMIT connlimit
| [ ENCRYPTED ] PASSWORD 'password' | PASSWORD NULL
| VALID UNTIL 'timestamp'
| PROFILE profile_name
| ACCOUNT UNLOCK | ACCOUNT LOCK
```

```
ALTER ROLE name RENAME TO new_name
```

```
ALTER ROLE { role_specification | ALL } [ IN DATABASE database_name ]
SET configuration_parameter { TO | = } { value | DEFAULT }
ALTER ROLE { role_specification | ALL } [ IN DATABASE database_name ]
SET configuration_parameter FROM CURRENT
ALTER ROLE { role_specification | ALL } [ IN DATABASE database_name ]
RESET configuration_parameter
ALTER ROLE { role_specification | ALL } [ IN DATABASE database_name ] RESET ALL
```

where *role_specification* can be:

```
role_name
| CURRENT_ROLE
| CURRENT_USER
| SESSION_USER
```

Description

ALTER ROLE changes the attributes of a Postgres Pro role.

The first variant of this command listed in the synopsis can change many of the role attributes that can be specified in [CREATE ROLE](#). (All the possible attributes are covered, except that there are no options for adding or removing memberships; use [GRANT](#) and [REVOKE](#) for that.) Attributes not mentioned in the command retain their previous settings. Database superusers can change any of these settings for any role, except for changing the SUPERUSER property for the [bootstrap superuser](#). Non-superuser roles having CREATEROLE privilege can change most of these properties, but only for non-superuser and non-replication roles for which they have been granted ADMIN OPTION. Non-superusers cannot change the SUPERUSER property and can change the CREATEDB, REPLICATION, and BYPASSRLS properties only if they possess the corresponding property themselves. Ordinary roles can only change their own password.

The second variant changes the name of the role. Database superusers can rename any role. Roles having CREATEROLE privilege can rename non-superuser roles for which they have been granted ADMIN OPTION. The current session user cannot be renamed. (Connect as a different user if you need to do that.)

Because MD5-encrypted passwords use the role name as cryptographic salt, renaming a role clears its password if the password is MD5-encrypted.

The remaining variants change a role's session default for a configuration variable, either for all databases or, when the `IN DATABASE` clause is specified, only for sessions in the named database. If `ALL` is specified instead of a role name, this changes the setting for all roles. Using `ALL` with `IN DATABASE` is effectively the same as using the command `ALTER DATABASE ... SET ...`.

Whenever the role subsequently starts a new session, the specified value becomes the session default, overriding whatever setting is present in `postgresql.conf` or has been received from the `postgres` command line. This only happens at login time; executing `SET ROLE` or `SET SESSION AUTHORIZATION` does not cause new configuration values to be set. Settings set for all databases are overridden by database-specific settings attached to a role. Settings for specific databases or specific roles override settings for all roles.

Superusers can change anyone's session defaults. Roles having `CREATEROLE` privilege can change defaults for non-superuser roles for which they have been granted `ADMIN OPTION`. Ordinary roles can only set defaults for themselves. Certain configuration variables cannot be set this way, or can only be set if a superuser issues the command. Only superusers can change a setting for all roles in all databases.

Parameters

name

The name of the role whose attributes are to be altered.

`CURRENT_ROLE`

`CURRENT_USER`

Alter the current user instead of an explicitly identified role.

`SESSION_USER`

Alter the current session user instead of an explicitly identified role.

`SUPERUSER`

`NOSUPERUSER`

`CREATEDB`

`NOCREATEDB`

`CREATEROLE`

`NOCREATEROLE`

`INHERIT`

`NOINHERIT`

`LOGIN`

`NOLOGIN`

`REPLICATION`

`NOREPLICATION`

`BYPASSRLS`

`NOBYPASSRLS`

`CONNECTION LIMIT` *conlimit*

`[ENCRYPTED] PASSWORD` '*password*'

`PASSWORD NULL`

`VALID UNTIL` '*timestamp*'

`PROFILE` *profile_name*

These clauses alter attributes originally set by `CREATE ROLE`. For more information, see the `CREATE ROLE` reference page.

`ACCOUNT UNLOCK`

Unlock the role. A role can be locked if the allowed number of failed login attempts has been exceeded (see `CREATE PROFILE` for details).

ACCOUNT LOCK

Explicitly lock the role.

new_name

The new name of the role.

database_name

The name of the database the configuration variable should be set in.

configuration_parameter
value

Set this role's session default for the specified configuration parameter to the given value. If *value* is `DEFAULT` or, equivalently, `RESET` is used, the role-specific variable setting is removed, so the role will inherit the system-wide default setting in new sessions. Use `RESET ALL` to clear all role-specific settings. `SET FROM CURRENT` saves the session's current value of the parameter as the role-specific value. If `IN DATABASE` is specified, the configuration parameter is set or removed for the given role and database only.

Role-specific variable settings take effect only at login; `SET ROLE` and `SET SESSION AUTHORIZATION` do not process role-specific variable settings.

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

Notes

Use [CREATE ROLE](#) to add new roles, and [DROP ROLE](#) to remove a role.

`ALTER ROLE` cannot change a role's memberships. Use [GRANT](#) and [REVOKE](#) to do that.

Caution must be exercised when specifying an unencrypted password with this command. The password will be transmitted to the server in cleartext, and it might also be logged in the client's command history or the server log. [psql](#) contains a command `\password` that can be used to change a role's password without exposing the cleartext password.

It is also possible to tie a session default to a specific database rather than to a role; see [ALTER DATABASE](#). If there is a conflict, database-role-specific settings override role-specific ones, which in turn override database-specific ones.

Examples

Change a role's password:

```
ALTER ROLE davide WITH PASSWORD 'hu8jmn3';
```

Remove a role's password:

```
ALTER ROLE davide WITH PASSWORD NULL;
```

Change a password expiration date, specifying that the password should expire at midday on 4th May 2015 using the time zone which is one hour ahead of UTC:

```
ALTER ROLE chris VALID UNTIL 'May 4 12:00:00 2015 +1';
```

Make a password valid forever:

```
ALTER ROLE fred VALID UNTIL 'infinity';
```

Give a role the ability to manage other roles and create new databases:

```
ALTER ROLE miriam CREATEROLE CREATEDB;
```

Give a role a non-default setting of the [maintenance_work_mem](#) parameter:

```
ALTER ROLE worker_bee SET maintenance_work_mem = 100000;
```

Give a role a non-default, database-specific setting of the [client_min_messages](#) parameter:

```
ALTER ROLE fred IN DATABASE devel SET client_min_messages = DEBUG;
```

Compatibility

The `ALTER ROLE` statement is a Postgres Pro extension.

See Also

[CREATE ROLE](#), [DROP ROLE](#), [ALTER DATABASE](#), [SET](#), [CREATE PROFILE](#)

ALTER ROUTINE

ALTER ROUTINE — change the definition of a routine

Synopsis

```
ALTER ROUTINE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]
    action [ ... ] [ RESTRICT ]
ALTER ROUTINE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]
    RENAME TO new_name
ALTER ROUTINE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
ALTER ROUTINE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]
    SET SCHEMA new_schema
ALTER ROUTINE name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]
    [ NO ] DEPENDS ON EXTENSION extension_name
```

where *action* is one of:

```
IMMUTABLE | STABLE | VOLATILE
[ NOT ] LEAKPROOF
[ EXTERNAL ] SECURITY INVOKER | [ EXTERNAL ] SECURITY DEFINER
PARALLEL { UNSAFE | RESTRICTED | SAFE }
COST execution_cost
ROWS result_rows
SET configuration_parameter { TO | = } { value | DEFAULT }
SET configuration_parameter FROM CURRENT
RESET configuration_parameter
RESET ALL
```

Description

ALTER ROUTINE changes the definition of a routine, which can be an aggregate function, a normal function, or a procedure. See under [ALTER AGGREGATE](#), [ALTER FUNCTION](#), and [ALTER PROCEDURE](#) for the description of the parameters, more examples, and further details.

Examples

To rename the routine `foo` for type `integer` to `foobar`:

```
ALTER ROUTINE foo(integer) RENAME TO foobar;
```

This command will work independent of whether `foo` is an aggregate, function, or procedure.

Compatibility

This statement is partially compatible with the ALTER ROUTINE statement in the SQL standard. See under [ALTER FUNCTION](#) and [ALTER PROCEDURE](#) for more details. Allowing routine names to refer to aggregate functions is a Postgres Pro extension.

See Also

[ALTER AGGREGATE](#), [ALTER FUNCTION](#), [ALTER PROCEDURE](#), [DROP ROUTINE](#)

Note that there is no CREATE ROUTINE command.

ALTER RULE

ALTER RULE — change the definition of a rule

Synopsis

```
ALTER RULE name ON table_name RENAME TO new_name
```

Description

ALTER RULE changes properties of an existing rule. Currently, the only available action is to change the rule's name.

To use ALTER RULE, you must own the table or view that the rule applies to.

Parameters

name

The name of an existing rule to alter.

table_name

The name (optionally schema-qualified) of the table or view that the rule applies to.

new_name

The new name for the rule.

Examples

To rename an existing rule:

```
ALTER RULE notify_all ON emp RENAME TO notify_me;
```

Compatibility

ALTER RULE is a Postgres Pro language extension, as is the entire query rewrite system.

See Also

[CREATE RULE](#), [DROP RULE](#)

ALTER SCHEMA

ALTER SCHEMA — change the definition of a schema

Synopsis

```
ALTER SCHEMA name RENAME TO new_name
ALTER SCHEMA name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
ALTER SCHEMA name SECURITY OFFICER TO new_security_officer
ALTER SCHEMA name RESET SECURITY OFFICER
```

Description

ALTER SCHEMA changes the definition of a schema.

You must own the schema to use ALTER SCHEMA. To rename a schema you must also have the CREATE privilege for the database. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have the CREATE privilege for the database. (Note that superusers have all these privileges automatically.)

The SECURITY OFFICER clause sets the security officer of the schema, while the RESET SECURITY OFFICER clause resets the security officer. When you set the security officer, the schema becomes the [vault](#) schema.

Parameters

name

The name of an existing schema.

new_name

The new name of the schema. The new name cannot begin with `pg_`, as such names are reserved for system schemas.

new_owner

The new owner of the schema.

new_security_officer

The new security officer of the schema.

Compatibility

There is no ALTER SCHEMA statement in the SQL standard.

See Also

[CREATE SCHEMA](#), [DROP SCHEMA](#)

ALTER SEQUENCE

ALTER SEQUENCE — change the definition of a sequence generator

Synopsis

For sharded objects:

```
ALTER SEQUENCE [ IF EXISTS ] name
    [ AS data_type ]
    [ INCREMENT [ BY ] increment ]
    [ MINVALUE minvalue | NO MINVALUE ] [ MAXVALUE maxvalue | NO MAXVALUE ]
    [ RESTART [ [ WITH ] restart ] ]
    [ CACHE cache ] [ [ NO ] CYCLE ]
    [ OWNED BY { table_name.column_name | NONE } ]
ALTER SEQUENCE [ IF EXISTS ] name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER SEQUENCE [ IF EXISTS ] name RENAME TO new_name
ALTER SEQUENCE [ IF EXISTS ] name SET SCHEMA new_schema
```

For other objects types:

```
ALTER SEQUENCE [ IF EXISTS ] name
    [ AS data_type ]
    [ INCREMENT [ BY ] increment ]
    [ MINVALUE minvalue | NO MINVALUE ] [ MAXVALUE maxvalue | NO MAXVALUE ]
    [ START [ WITH ] start ]
    [ RESTART [ [ WITH ] restart ] ]
    [ CACHE cache ] [ [ NO ] CYCLE ]
    [ OWNED BY { table_name.column_name | NONE } ]
ALTER SEQUENCE [ IF EXISTS ] name SET { LOGGED | UNLOGGED }
ALTER SEQUENCE [ IF EXISTS ] name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER SEQUENCE [ IF EXISTS ] name RENAME TO new_name
ALTER SEQUENCE [ IF EXISTS ] name SET SCHEMA new_schema
```

Description

ALTER SEQUENCE changes the parameters of an existing sequence generator. Any parameters not specifically set in the ALTER SEQUENCE command retain their prior settings.

You must own the sequence to use ALTER SEQUENCE. To change a sequence's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the sequence's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the sequence. However, a superuser can alter ownership of any sequence anyway.)

Parameters

name

The name (optionally schema-qualified) of a sequence to be altered.

IF EXISTS

Do not throw an error if the sequence does not exist. A notice is issued in this case.

data_type

The optional clause AS *data_type* changes the data type of the sequence. Valid types are smallint, integer, and bigint.

Changing the data type automatically changes the minimum and maximum values of the sequence if and only if the previous minimum and maximum values were the minimum or maximum value of the old data type (in other words, if the sequence had been created using `NO MINVALUE` or `NO MAXVALUE`, implicitly or explicitly). Otherwise, the minimum and maximum values are preserved, unless new values are given as part of the same command. If the minimum and maximum values do not fit into the new data type, an error will be generated.

increment

The clause `INCREMENT BY increment` is optional. A positive value will make an ascending sequence, a negative one a descending sequence. If unspecified, the old increment value will be maintained.

minvalue

`NO MINVALUE`

The optional clause `MINVALUE minvalue` determines the minimum value a sequence can generate. If `NO MINVALUE` is specified, the defaults of 1 and the minimum value of the data type for ascending and descending sequences, respectively, will be used. If neither option is specified, the current minimum value will be maintained.

maxvalue

`NO MAXVALUE`

The optional clause `MAXVALUE maxvalue` determines the maximum value for the sequence. If `NO MAXVALUE` is specified, the defaults of the maximum value of the data type and -1 for ascending and descending sequences, respectively, will be used. If neither option is specified, the current maximum value will be maintained.

start

The optional clause `START WITH start` changes the recorded start value of the sequence. This has no effect on the *current* sequence value; it simply sets the value that future `ALTER SEQUENCE RESTART` commands will use.

restart

The optional clause `RESTART [WITH restart]` changes the current value of the sequence. This is similar to calling the `setval` function with `is_called = false`: the specified value will be returned by the *next* call of `nextval`. Writing `RESTART` with no *restart* value is equivalent to supplying the start value that was recorded by `CREATE SEQUENCE` or last set by `ALTER SEQUENCE START WITH`.

In contrast to a `setval` call, a `RESTART` operation on a sequence is transactional and blocks concurrent transactions from obtaining numbers from the same sequence. If that's not the desired mode of operation, `setval` should be used.

cache

The clause `CACHE cache` enables sequence numbers to be preallocated and stored in memory for faster access. The minimum value is 1 (only one value can be generated at a time, i.e., no cache). If unspecified, the old cache value will be maintained.

`CYCLE`

The optional `CYCLE` key word can be used to enable the sequence to wrap around when the *maxvalue* or *minvalue* has been reached by an ascending or descending sequence respectively. If the limit is reached, the next number generated will be the *minvalue* or *maxvalue*, respectively.

`NO CYCLE`

If the optional `NO CYCLE` key word is specified, any calls to `nextval` after the sequence has reached its maximum value will return an error. If neither `CYCLE` or `NO CYCLE` are specified, the old cycle behavior will be maintained.

SET { LOGGED | UNLOGGED }

This form changes the sequence from unlogged to logged or vice-versa (see [CREATE SEQUENCE](#)). It cannot be applied to a temporary sequence.

OWNED BY *table_name.column_name*

OWNED BY NONE

The OWNED BY option causes the sequence to be associated with a specific table column, such that if that column (or its whole table) is dropped, the sequence will be automatically dropped as well. If specified, this association replaces any previously specified association for the sequence. The specified table must have the same owner and be in the same schema as the sequence. Specifying OWNED BY NONE removes any existing association, making the sequence “free-standing”.

new_owner

The user name of the new owner of the sequence.

new_name

The new name for the sequence.

new_schema

The new schema for the sequence.

RESTART WITH

The minimum sequence value parameter in Postgres Pro Shardman works more like a lower boundary on the global interval of available values, so it can only be increased to make sure no duplicate numbers are generated.

The RESTART WITH clause allows restarting a sequence at any arbitrary lower bound, but in this case, there is no guarantee that previously generated numbers will not repeat.

Using both RESTART WITH and MINVALUE in a single statement is not permitted to avoid confusion.

Notes

ALTER SEQUENCE will not immediately affect `nextval` results in backends, other than the current one, that have preallocated (cached) sequence values. They will use up all cached values prior to noticing the changed sequence generation parameters. The current backend will be affected immediately.

ALTER SEQUENCE does not affect the `currval` status for the sequence. (Before PostgreSQL 8.3, it sometimes did.)

ALTER SEQUENCE blocks concurrent `nextval`, `currval`, `lastval`, and `setval` calls.

For historical reasons, ALTER TABLE can be used with sequences too; but the only variants of ALTER TABLE that are allowed with sequences are equivalent to the forms shown above.

Examples

Alter the block size parameter of a sequence called `serial`:

```
ALTER SEQUENCE serial SET (block_size = 8192);
```

Restart a sequence called `serial`, at 105:

```
ALTER SEQUENCE serial RESTART WITH 105;
```

Compatibility

ALTER SEQUENCE conforms to the SQL standard, except for the AS, START WITH, OWNED BY, OWNER TO, RENAME TO, and SET SCHEMA clauses, which are Postgres Pro extensions.

See Also

[CREATE SEQUENCE](#), [DROP SEQUENCE](#)

ALTER SERVER

ALTER SERVER — change the definition of a foreign server

Synopsis

```
ALTER SERVER name [ VERSION 'new_version' ]  
    [ OPTIONS ( [ ADD | SET | DROP ] option ['value'] [, ... ] ) ]  
ALTER SERVER name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }  
ALTER SERVER name RENAME TO new_name
```

Description

ALTER SERVER changes the definition of a foreign server. The first form changes the server version string or the generic options of the server (at least one clause is required). The second form changes the owner of the server.

To alter the server you must be the owner of the server. Additionally to alter the owner, you must be able to SET ROLE to the new owning role, and you must have USAGE privilege on the server's foreign-data wrapper. (Note that superusers satisfy all these criteria automatically.)

Parameters

name

The name of an existing server.

new_version

New server version.

OPTIONS ([ADD | SET | DROP] *option* ['*value*'] [, ...])

Change options for the server. ADD, SET, and DROP specify the action to be performed. ADD is assumed if no operation is explicitly specified. Option names must be unique; names and values are also validated using the server's foreign-data wrapper library.

new_owner

The user name of the new owner of the foreign server.

new_name

The new name for the foreign server.

Examples

Alter server `foo`, add connection options:

```
ALTER SERVER foo OPTIONS (host 'foo', dbname 'foodb');
```

Alter server `foo`, change version, change host option:

```
ALTER SERVER foo VERSION '8.4' OPTIONS (SET host 'baz');
```

Compatibility

ALTER SERVER conforms to ISO/IEC 9075-9 (SQL/MED). The OWNER TO and RENAME forms are Postgres Pro extensions.

See Also

[CREATE SERVER](#), [DROP SERVER](#)

ALTER STATISTICS

ALTER STATISTICS — change the definition of an extended statistics object

Synopsis

```
ALTER STATISTICS name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |  
SESSION_USER }  
ALTER STATISTICS name RENAME TO new_name  
ALTER STATISTICS name SET SCHEMA new_schema  
ALTER STATISTICS name SET STATISTICS { new_target | DEFAULT }
```

Description

ALTER STATISTICS changes the parameters of an existing extended statistics object. Any parameters not specifically set in the ALTER STATISTICS command retain their prior settings.

You must own the statistics object to use ALTER STATISTICS. To change a statistics object's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the statistics object's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the statistics object. However, a superuser can alter ownership of any statistics object anyway.)

Parameters

name

The name (optionally schema-qualified) of the statistics object to be altered.

new_owner

The user name of the new owner of the statistics object.

new_name

The new name for the statistics object.

new_schema

The new schema for the statistics object.

new_target

The statistic-gathering target for this statistics object for subsequent [ANALYZE](#) operations. The target can be set in the range 0 to 10000. Set it to DEFAULT to revert to using the system default statistics target ([default_statistics_target](#)). (Setting to a value of -1 is an obsolete way spelling to get the same outcome.) For more information on the use of statistics by the Postgres Pro query planner, refer to [Section 14.2](#).

Compatibility

There is no ALTER STATISTICS command in the SQL standard.

See Also

[CREATE STATISTICS](#), [DROP STATISTICS](#)

ALTER SUBSCRIPTION

ALTER SUBSCRIPTION — change the definition of a subscription

Synopsis

```
ALTER SUBSCRIPTION name CONNECTION 'conninfo'
ALTER SUBSCRIPTION name SET PUBLICATION publication_name [, ...] [ WITH
( publication_option [= value] [, ... ] ) ]
ALTER SUBSCRIPTION name ADD PUBLICATION publication_name [, ...] [ WITH
( publication_option [= value] [, ... ] ) ]
ALTER SUBSCRIPTION name DROP PUBLICATION publication_name [, ...] [ WITH
( publication_option [= value] [, ... ] ) ]
ALTER SUBSCRIPTION name REFRESH PUBLICATION [ WITH ( refresh_option [= value]
[, ... ] ) ]
ALTER SUBSCRIPTION name ENABLE
ALTER SUBSCRIPTION name DISABLE
ALTER SUBSCRIPTION name SET ( subscription_parameter [= value] [, ... ] )
ALTER SUBSCRIPTION name SKIP ( skip_option = value )
ALTER SUBSCRIPTION name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER SUBSCRIPTION name RENAME TO new_name
```

Description

ALTER SUBSCRIPTION can change most of the subscription properties that can be specified in [CREATE SUBSCRIPTION](#).

You must own the subscription to use ALTER SUBSCRIPTION. To rename a subscription or alter the owner, you must have CREATE permission on the database. In addition, to alter the owner, you must be able to SET ROLE to the new owning role. If the subscription has password_required=false, only superusers can modify it.

When refreshing a publication we remove the relations that are no longer part of the publication and we also remove the table synchronization slots if there are any. It is necessary to remove these slots so that the resources allocated for the subscription on the remote host are released. If due to network breakdown or some other error, Postgres Pro is unable to remove the slots, an error will be reported. To proceed in this situation, the user either needs to retry the operation or disassociate the slot from the subscription and drop the subscription as explained in [DROP SUBSCRIPTION](#).

Commands ALTER SUBSCRIPTION ... REFRESH PUBLICATION, ALTER SUBSCRIPTION ... {SET|ADD|DROP} PUBLICATION ... with refresh option as true and ALTER SUBSCRIPTION ... SET (failover = true|false) cannot be executed inside a transaction block.

Commands ALTER SUBSCRIPTION ... REFRESH PUBLICATION and ALTER SUBSCRIPTION ... {SET|ADD|DROP} PUBLICATION ... with refresh option as true also cannot be executed when the subscription has two_phase commit enabled, unless copy_data is false. See column subtwophasestate of pg_subscription to know the actual two-phase state.

Parameters

name

The name of a subscription whose properties are to be altered.

CONNECTION '*conninfo*'

This clause replaces the connection string originally set by [CREATE SUBSCRIPTION](#). See there for more information.

```
SET PUBLICATION publication_name
ADD PUBLICATION publication_name
DROP PUBLICATION publication_name
```

These forms change the list of subscribed publications. `SET` replaces the entire list of publications with a new list, `ADD` adds additional publications to the list of publications, and `DROP` removes the publications from the list of publications. We allow non-existent publications to be specified in `ADD` and `SET` variants so that users can add those later. See [CREATE SUBSCRIPTION](#) for more information. By default, this command will also act like `REFRESH PUBLICATION`.

publication_option specifies additional options for this operation. The supported options are:

`refresh` (boolean)

When false, the command will not try to refresh table information. `REFRESH PUBLICATION` should then be executed separately. The default is `true`.

Additionally, the options described under `REFRESH PUBLICATION` may be specified, to control the implicit refresh operation.

`REFRESH PUBLICATION`

Fetch missing table information from publisher. This will start replication of tables that were added to the subscribed-to publications since [CREATE SUBSCRIPTION](#) or the last invocation of `REFRESH PUBLICATION`.

refresh_option specifies additional options for the refresh operation. The supported options are:

`copy_data` (boolean)

Specifies whether to copy pre-existing data in the publications that are being subscribed to when the replication starts. The default is `true`.

Previously subscribed tables are not copied, even if a table's row filter `WHERE` clause has since been modified.

See [Notes](#) for details of how `copy_data = true` can interact with the `origin` parameter.

See the `binary` parameter of `CREATE SUBSCRIPTION` for details about copying pre-existing data in binary format.

`ENABLE`

Enables a previously disabled subscription, starting the logical replication worker at the end of the transaction.

`DISABLE`

Disables a running subscription, stopping the logical replication worker at the end of the transaction.

```
SET ( subscription_parameter [= value] [, ... ] )
```

This clause alters parameters originally set by [CREATE SUBSCRIPTION](#). See there for more information. The parameters that can be altered are `slot_name`, `synchronous_commit`, `binary`, `streaming`, `disable_on_error`, `password_required`, `run_as_owner`, `origin`, and `failover`. Only a superuser can set `password_required = false`.

When altering the `slot_name`, the `failover` and `two_phase` property values of the named slot may differ from the counterpart `failover` and `two_phase` parameters specified in the subscription. When creating the slot, ensure the slot properties `failover` and `two_phase` match their counterpart parameters of the subscription. Otherwise, the slot on the publisher may behave differently from what these subscription options say: for example, the slot on the publisher could either be synced to the standbys even when the subscription's `failover` option is disabled or could be disabled for sync even when the subscription's `failover` option is enabled.

`SKIP (skip_option = value)`

Skips applying all changes of the remote transaction. If incoming data violates any constraints, logical replication will stop until it is resolved. By using the `ALTER SUBSCRIPTION ... SKIP` command, the logical replication worker skips all data modification changes within the transaction. This option has no effect on the transactions that are already prepared by enabling [two_phase](#) on the subscriber. After the logical replication worker successfully skips the transaction or finishes a transaction, the LSN (stored in `pg_subscription.subskiplsn`) is cleared. See [Section 29.6](#) for the details of logical replication conflicts.

skip_option specifies options for this operation. The supported option is:

`lsn (pg_lsn)`

Specifies the finish LSN of the remote transaction whose changes are to be skipped by the logical replication worker. The finish LSN is the LSN at which the transaction is either committed or prepared. Skipping individual subtransactions is not supported. Setting `NONE` resets the LSN.

new_owner

The user name of the new owner of the subscription.

new_name

The new name for the subscription.

When specifying a parameter of type `boolean`, the `= value` part can be omitted, which is equivalent to specifying `TRUE`.

Examples

Change the publication subscribed by a subscription to `insert_only`:

```
ALTER SUBSCRIPTION mysub SET PUBLICATION insert_only;
```

Disable (stop) the subscription:

```
ALTER SUBSCRIPTION mysub DISABLE;
```

Compatibility

`ALTER SUBSCRIPTION` is a Postgres Pro extension.

See Also

[CREATE SUBSCRIPTION](#), [DROP SUBSCRIPTION](#), [CREATE PUBLICATION](#), [ALTER PUBLICATION](#)

ALTER TABLE

ALTER TABLE — change the definition of a table

Synopsis

For sharded tables:

```
ALTER TABLE [ IF EXISTS ] name [ * ]
    action [, ... ]
ALTER TABLE [ IF EXISTS ] name [ * ]
    RENAME [ COLUMN ] column_name TO new_column_name
ALTER TABLE [ IF EXISTS ] name [ * ]
    RENAME CONSTRAINT constraint_name TO new_constraint_name
ALTER TABLE [ IF EXISTS ] name
    RENAME TO new_name
ALTER TABLE [ IF EXISTS ] name
    SET SCHEMA new_schema
```

where *action* is one of:

```
    ADD [ COLUMN ] [ IF NOT EXISTS ] column_name data_type [ COLLATE collation ]
[ column_constraint [ ... ] ]
    DROP [ COLUMN ] [ IF EXISTS ] column_name [ RESTRICT | CASCADE ]
    ALTER [ COLUMN ] column_name [ SET DATA ] TYPE data_type [ COLLATE collation ]
[ USING expression ]
    ALTER [ COLUMN ] column_name SET DEFAULT expression
    ALTER [ COLUMN ] column_name DROP DEFAULT
    ALTER [ COLUMN ] column_name { SET | DROP } NOT NULL
    ALTER [ COLUMN ] column_name SET STATISTICS integer

    ADD table_constraint [ NOT VALID ]
    ADD table_constraint_using_index
    VALIDATE CONSTRAINT constraint_name
    DROP CONSTRAINT [ IF EXISTS ] constraint_name [ RESTRICT | CASCADE ]
    SET TABLESPACE new_tablespace
    RESET ( storage_parameter [, ... ] )
    OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

and *column_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ NOT NULL |
  NULL |
  CHECK ( expression ) [ NO INHERIT ] |
  DEFAULT default_expr |
  UNIQUE index_parameters |
  PRIMARY KEY index_parameters |
  REFERENCES reftable [ ( refcolumn ) ] [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ]
  [ ON DELETE referential_action ] [ ON UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *sharded_table_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ CHECK ( expression ) [ NO INHERIT ] |
  UNIQUE ( column_name [, ... ] ) index_parameters |
```


ALTER TABLE

```
EXCLUDE [ USING index_method ] ( exclude_element WITH operator
[, ... ] ) index_parameters [ WHERE ( predicate ) ] |
FOREIGN KEY ( column_name [, ... ] ) REFERENCES reftable [ ( refcolumn [, ... ] ) ]
[ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ] [ ON DELETE referential_action ] [ ON
UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *sharded_table_constraint_using_index* is:

```
[ CONSTRAINT constraint_name ]
{ UNIQUE | PRIMARY KEY } USING INDEX index_name
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

index_parameters in UNIQUE, PRIMARY KEY, and EXCLUDE constraints are:

```
[ INCLUDE ( column_name [, ... ] ) ]
[ WITH ( storage_parameter [= value] [, ... ] ) ]
[ USING INDEX TABLESPACE tablespace_name ]
```

exclude_element in an EXCLUDE constraint is:

```
{ column_name | ( expression ) } [ opclass ] [ ASC | DESC ] [ NULLS { FIRST | LAST } ]
```

referential_action in a FOREIGN KEY/REFERENCES constraint is:

```
{ NO ACTION | RESTRICT | CASCADE | SET NULL [ ( column_name [, ... ] ) ] | SET DEFAULT
[ ( column_name [, ... ] ) ] }
```

For other table types:

```
ALTER TABLE [ IF EXISTS ] [ ONLY ] name [ * ]
    action [, ... ]
ALTER TABLE [ IF EXISTS ] [ ONLY ] name [ * ]
    RENAME [ COLUMN ] column_name TO new_column_name
ALTER TABLE [ IF EXISTS ] [ ONLY ] name [ * ]
    RENAME CONSTRAINT constraint_name TO new_constraint_name
ALTER TABLE [ IF EXISTS ] name
    RENAME TO new_name
ALTER TABLE [ IF EXISTS ] name
    SET SCHEMA new_schema
ALTER TABLE ALL IN TABLESPACE name [ OWNED BY role_name [, ... ] ]
    SET TABLESPACE new_tablespace [ NOWAIT ]
ALTER TABLE [ IF EXISTS ] name
    ATTACH PARTITION partition_name { FOR VALUES partition_bound_spec | DEFAULT }
ALTER TABLE [ IF EXISTS ] name
    DETACH PARTITION partition_name [ CONCURRENTLY | FINALIZE ]
ALTER TABLE [ IF EXISTS ] name
    SPLIT PARTITION partition_name INTO
        (PARTITION partition_name1 { FOR VALUES partition_bound_spec | DEFAULT },
        PARTITION partition_name2 { FOR VALUES partition_bound_spec | DEFAULT }
    [, ...])
ALTER TABLE [ IF EXISTS ] name
    MERGE PARTITIONS (partition_name1, partition_name2 [, ...])
    INTO partition_name
```

where *action* is one of:

ALTER TABLE

```
ADD [ COLUMN ] [ IF NOT EXISTS ] column_name data_type [ COLLATE collation ]
[ column_constraint [ ... ] ]
DROP [ COLUMN ] [ IF EXISTS ] column_name [ RESTRICT | CASCADE ]
ALTER [ COLUMN ] column_name [ SET DATA ] TYPE data_type [ COLLATE collation ]
[ USING expression ]
ALTER [ COLUMN ] column_name SET DEFAULT expression
ALTER [ COLUMN ] column_name DROP DEFAULT
ALTER [ COLUMN ] column_name { SET | DROP } NOT NULL
ALTER [ COLUMN ] column_name SET EXPRESSION AS ( expression )
ALTER [ COLUMN ] column_name DROP EXPRESSION [ IF EXISTS ]
ALTER [ COLUMN ] column_name ADD GENERATED { ALWAYS | BY DEFAULT } AS IDENTITY
[ ( sequence_options ) ]
ALTER [ COLUMN ] column_name { SET GENERATED { ALWAYS | BY DEFAULT } |
SET sequence_option | RESTART [ [ WITH ] restart ] } [...]
ALTER [ COLUMN ] column_name DROP IDENTITY [ IF EXISTS ]
ALTER [ COLUMN ] column_name SET STATISTICS { integer | DEFAULT }
ALTER [ COLUMN ] column_name SET ( attribute_option = value [, ... ] )
ALTER [ COLUMN ] column_name RESET ( attribute_option [, ... ] )
ALTER [ COLUMN ] column_name SET STORAGE { PLAIN | EXTERNAL | EXTENDED | MAIN |
DEFAULT }
ALTER [ COLUMN ] column_name SET COMPRESSION compression_method
ADD table_constraint [ NOT VALID ]
ADD table_constraint_using_index
ALTER CONSTRAINT constraint_name [ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY
DEFERRED | INITIALLY IMMEDIATE ]
VALIDATE CONSTRAINT constraint_name
DROP CONSTRAINT [ IF EXISTS ] constraint_name [ RESTRICT | CASCADE ]
DISABLE TRIGGER [ trigger_name | ALL | USER ]
ENABLE TRIGGER [ trigger_name | ALL | USER ]
ENABLE REPLICA TRIGGER trigger_name
ENABLE ALWAYS TRIGGER trigger_name
DISABLE RULE rewrite_rule_name
ENABLE RULE rewrite_rule_name
ENABLE REPLICA RULE rewrite_rule_name
ENABLE ALWAYS RULE rewrite_rule_name
DISABLE ROW LEVEL SECURITY
ENABLE ROW LEVEL SECURITY
FORCE ROW LEVEL SECURITY
NO FORCE ROW LEVEL SECURITY
CLUSTER ON index_name
SET WITHOUT CLUSTER
SET WITHOUT OIDS
SET ACCESS METHOD { new_access_method | DEFAULT }
SET TABLESPACE new_tablespace
SET { LOGGED | UNLOGGED }
SET CONSTANT
SET ( storage_parameter [= value] [, ... ] )
RESET ( storage_parameter [, ... ] )
INHERIT parent_table
NO INHERIT parent_table
OF type_name
NOT OF
OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
REPLICA IDENTITY { DEFAULT | USING INDEX index_name | FULL | NOTHING }
```

and *partition_bound_spec* is:

```
IN ( partition_bound_expr [, ...] ) |
```

```
FROM ( { partition_bound_expr | MINVALUE | MAXVALUE } [, ...] )
TO ( { partition_bound_expr | MINVALUE | MAXVALUE } [, ...] ) |
WITH ( MODULUS numeric_literal, REMAINDER numeric_literal )
```

and *column_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ NOT NULL |
  NULL |
  CHECK ( expression ) [ NO INHERIT ] |
  DEFAULT default_expr |
  GENERATED ALWAYS AS ( generation_expr ) STORED |
  GENERATED { ALWAYS | BY DEFAULT } AS IDENTITY [ ( sequence_options ) ] |
  UNIQUE [ NULLS [ NOT ] DISTINCT ] index_parameters |
  PRIMARY KEY index_parameters |
  REFERENCES reftable [ ( refcolumn ) ] [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ]
  [ ON DELETE referential_action ] [ ON UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *table_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ CHECK ( expression ) [ NO INHERIT ] |
  UNIQUE [ NULLS [ NOT ] DISTINCT ] ( column_name [, ...] ) index_parameters |
  PRIMARY KEY ( column_name [, ...] ) index_parameters |
  EXCLUDE [ USING index_method ] ( exclude_element WITH operator
  [, ...] ) index_parameters [ WHERE ( predicate ) ] |
  FOREIGN KEY ( column_name [, ...] ) REFERENCES reftable [ ( refcolumn [, ...] ) ]
  [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ] [ ON DELETE referential_action ] [ ON
  UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *table_constraint_using_index* is:

```
[ CONSTRAINT constraint_name ]
{ UNIQUE | PRIMARY KEY } USING INDEX index_name
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

index_parameters in UNIQUE, PRIMARY KEY, and EXCLUDE constraints are:

```
[ INCLUDE ( column_name [, ...] ) ]
[ WITH ( storage_parameter [= value] [, ...] ) ]
[ USING INDEX TABLESPACE tablespace_name ]
```

exclude_element in an EXCLUDE constraint is:

```
{ column_name | ( expression ) } [ COLLATE collation ] [ opclass [ ( opclass_parameter
= value [, ...] ) ] ] [ ASC | DESC ] [ NULLS { FIRST | LAST } ]
```

referential_action in a FOREIGN KEY/REFERENCES constraint is:

```
{ NO ACTION | RESTRICT | CASCADE | SET NULL [ ( column_name [, ...] ) ] | SET DEFAULT
[ ( column_name [, ...] ) ] }
```

Description

ALTER TABLE changes the definition of an existing table. There are several subforms described below. Note that the lock level required may differ for each subform. An ACCESS EXCLUSIVE lock is acquired

unless explicitly noted. When multiple subcommands are given, the lock acquired will be the strictest one required by any subcommand.

ADD COLUMN [IF NOT EXISTS]

This form adds a new column to the table, using the same syntax as [CREATE TABLE](#). If `IF NOT EXISTS` is specified and a column already exists with this name, no error is thrown.

DROP COLUMN [IF EXISTS]

This form drops a column from a table. Indexes and table constraints involving the column will be automatically dropped as well. Multivariate statistics referencing the dropped column will also be removed if the removal of the column would cause the statistics to contain data for only a single column. You will need to say `CASCADE` if anything outside the table depends on the column, for example, foreign key references or views. If `IF EXISTS` is specified and the column does not exist, no error is thrown. In this case a notice is issued instead.

SET DATA TYPE

This form changes the type of a column of a table. Indexes and simple table constraints involving the column will be automatically converted to use the new column type by reparsing the originally supplied expression. The optional `COLLATE` clause specifies a collation for the new column; if omitted, the collation is the default for the new column type. The optional `USING` clause specifies how to compute the new column value from the old; if omitted, the default conversion is the same as an assignment cast from old data type to new. A `USING` clause must be provided if there is no implicit or assignment cast from old to new type.

When this form is used, the column's statistics are removed, so running [ANALYZE](#) on the table afterwards is recommended.

SET/DROP DEFAULT

These forms set or remove the default value for a column (where removal is equivalent to setting the default value to `NULL`). The new default value will only apply in subsequent `INSERT` or `UPDATE` commands; it does not cause rows already in the table to change.

SET/DROP NOT NULL

These forms change whether a column is marked to allow null values or to reject null values.

`SET NOT NULL` may only be applied to a column provided none of the records in the table contain a `NULL` value for the column. Ordinarily this is checked during the `ALTER TABLE` by scanning the entire table; however, if a valid `CHECK` constraint is found which proves no `NULL` can exist, then the table scan is skipped.

If this table is a partition, one cannot perform `DROP NOT NULL` on a column if it is marked `NOT NULL` in the parent table. To drop the `NOT NULL` constraint from all the partitions, perform `DROP NOT NULL` on the parent table. Even if there is no `NOT NULL` constraint on the parent, such a constraint can still be added to individual partitions, if desired; that is, the children can disallow nulls even if the parent allows them, but not the other way around.

SET EXPRESSION AS

This form replaces the expression of a generated column. Existing data in the column is rewritten and all the future changes will apply the new generation expression.

DROP EXPRESSION [IF EXISTS]

This form turns a stored generated column into a normal base column. Existing data in the columns is retained, but future changes will no longer apply the generation expression.

If `DROP EXPRESSION IF EXISTS` is specified and the column is not a stored generated column, no error is thrown. In this case a notice is issued instead.

```
ADD GENERATED { ALWAYS | BY DEFAULT } AS IDENTITY
SET GENERATED { ALWAYS | BY DEFAULT }
DROP IDENTITY [ IF EXISTS ]
```

These forms change whether a column is an identity column or change the generation attribute of an existing identity column. See [CREATE TABLE](#) for details. Like `SET DEFAULT`, these forms only affect the behavior of subsequent `INSERT` and `UPDATE` commands; they do not cause rows already in the table to change.

If `DROP IDENTITY IF EXISTS` is specified and the column is not an identity column, no error is thrown. In this case a notice is issued instead.

```
SET sequence_option
RESTART
```

These forms alter the sequence that underlies an existing identity column. *sequence_option* is an option supported by [ALTER SEQUENCE](#) such as `INCREMENT BY`.

```
SET STATISTICS
```

This form sets the per-column statistics-gathering target for subsequent [ANALYZE](#) operations. The target can be set in the range 0 to 10000. Set it to `DEFAULT` to revert to using the system default statistics target ([default_statistics_target](#)). (Setting to a value of -1 is an obsolete way spelling to get the same outcome.) For more information on the use of statistics by the Postgres Pro query planner, refer to [Section 14.2](#).

`SET STATISTICS` acquires a `SHARE UPDATE EXCLUSIVE` lock.

```
SET ( attribute_option = value [, ... ] )
RESET ( attribute_option [, ... ] )
```

This form sets or resets per-attribute options. Currently, the only defined per-attribute options are `n_distinct` and `n_distinct_inherited`, which override the number-of-distinct-values estimates made by subsequent [ANALYZE](#) operations. `n_distinct` affects the statistics for the table itself, while `n_distinct_inherited` affects the statistics gathered for the table plus its inheritance children. When set to a positive value, `ANALYZE` will assume that the column contains exactly the specified number of distinct nonnull values. When set to a negative value, which must be greater than or equal to -1, `ANALYZE` will assume that the number of distinct nonnull values in the column is linear in the size of the table; the exact count is to be computed by multiplying the estimated table size by the absolute value of the given number. For example, a value of -1 implies that all values in the column are distinct, while a value of -0.5 implies that each value appears twice on the average. This can be useful when the size of the table changes over time, since the multiplication by the number of rows in the table is not performed until query planning time. Specify a value of 0 to revert to estimating the number of distinct values normally. For more information on the use of statistics by the Postgres Pro query planner, refer to [Section 14.2](#).

Changing per-attribute options acquires a `SHARE UPDATE EXCLUSIVE` lock.

```
SET STORAGE { PLAIN | EXTERNAL | EXTENDED | MAIN | DEFAULT }
```

This form sets the storage mode for a column. This controls whether this column is held inline or in a secondary TOAST table, and whether the data should be compressed or not. `PLAIN` must be used for fixed-length values such as `integer` and is inline, uncompressed. `MAIN` is for inline, compressible data. `EXTERNAL` is for external, uncompressed data, and `EXTENDED` is for external, compressed data. Writing `DEFAULT` sets the storage mode to the default mode for the column's data type. `EXTENDED` is the default for most data types that support non-`PLAIN` storage. Use of `EXTERNAL` will make substring operations on very large `text` and `bytea` values run faster, at the penalty of increased storage space. Note that `ALTER TABLE ... SET STORAGE` doesn't itself change anything in the table; it just sets the strategy to be pursued during future table updates. See [Section 66.2](#) for more information.

```
SET COMPRESSION compression_method
```

This form sets the compression method for a column, determining how values inserted in future will be compressed (if the storage mode permits compression at all). This does not cause the table to be rewritten, so existing data may still be compressed with other compression methods. If the table is restored with `pg_restore`, then all values are rewritten with the configured compression method. However, when data is inserted from another relation (for example, by `INSERT ... SELECT`), values from the source table are not necessarily detoasted, so any previously compressed data may retain its existing compression method, rather than being recompressed with the compression method of the target column. The supported compression methods are `pglz` and `lz4`. (`lz4` is available only if `--with-lz4` was used when building Postgres Pro.) In addition, *compression_method* can be `default`, which selects the default behavior of consulting the `default_toast_compression` setting at the time of data insertion to determine the method to use.

```
ADD table_constraint [ NOT VALID ]
```

This form adds a new constraint to a table using the same constraint syntax as `CREATE TABLE`, plus the option `NOT VALID`, which is currently only allowed for foreign key and CHECK constraints.

Normally, this form will cause a scan of the table to verify that all existing rows in the table satisfy the new constraint. But if the `NOT VALID` option is used, this potentially-lengthy scan is skipped. The constraint will still be enforced against subsequent inserts or updates (that is, they'll fail unless there is a matching row in the referenced table, in the case of foreign keys, or they'll fail unless the new row matches the specified check condition). But the database will not assume that the constraint holds for all rows in the table, until it is validated by using the `VALIDATE CONSTRAINT` option. See [Notes](#) below for more information about using the `NOT VALID` option.

Although most forms of `ADD table_constraint` require an `ACCESS EXCLUSIVE` lock, `ADD FOREIGN KEY` requires only a `SHARE ROW EXCLUSIVE` lock. Note that `ADD FOREIGN KEY` also acquires a `SHARE ROW EXCLUSIVE` lock on the referenced table, in addition to the lock on the table on which the constraint is declared.

Additional restrictions apply when unique or primary key constraints are added to partitioned tables; see `CREATE TABLE`. Also, foreign key constraints on partitioned tables may not be declared `NOT VALID` at present.

```
ADD table_constraint_using_index
```

This form adds a new `PRIMARY KEY` or `UNIQUE` constraint to a table based on an existing unique index. All the columns of the index will be included in the constraint.

The index cannot have expression columns nor be a partial index. Also, it must be a b-tree index with default sort ordering. These restrictions ensure that the index is equivalent to one that would be built by a regular `ADD PRIMARY KEY` or `ADD UNIQUE` command.

If `PRIMARY KEY` is specified, and the index's columns are not already marked `NOT NULL`, then this command will attempt to do `ALTER COLUMN SET NOT NULL` against each such column. That requires a full table scan to verify the column(s) contain no nulls. In all other cases, this is a fast operation.

If a constraint name is provided then the index will be renamed to match the constraint name. Otherwise the constraint will be named the same as the index.

After this command is executed, the index is “owned” by the constraint, in the same way as if the index had been built by a regular `ADD PRIMARY KEY` or `ADD UNIQUE` command. In particular, dropping the constraint will make the index disappear too.

This form is not currently supported on partitioned tables.

Note

Adding a constraint using an existing index can be helpful in situations where a new constraint needs to be added without blocking table updates for a long time. To do that, create the index

using `CREATE UNIQUE INDEX CONCURRENTLY`, and then convert it to a constraint using this syntax. See the example below.

ALTER CONSTRAINT

This form alters the attributes of a constraint that was previously created. Currently only foreign key constraints may be altered.

VALIDATE CONSTRAINT

This form validates a foreign key or check constraint that was previously created as `NOT VALID`, by scanning the table to ensure there are no rows for which the constraint is not satisfied. Nothing happens if the constraint is already marked valid. (See [Notes](#) below for an explanation of the usefulness of this command.)

This command acquires a `SHARE UPDATE EXCLUSIVE` lock.

DROP CONSTRAINT [IF EXISTS]

This form drops the specified constraint on a table, along with any index underlying the constraint. If `IF EXISTS` is specified and the constraint does not exist, no error is thrown. In this case a notice is issued instead.

DISABLE/ENABLE [REPLICA | ALWAYS] TRIGGER

These forms configure the firing of trigger(s) belonging to the table. A disabled trigger is still known to the system, but is not executed when its triggering event occurs. (For a deferred trigger, the enable status is checked when the event occurs, not when the trigger function is actually executed.) One can disable or enable a single trigger specified by name, or all triggers on the table, or only user triggers (this option excludes internally generated constraint triggers, such as those that are used to implement foreign key constraints or deferrable uniqueness and exclusion constraints). Disabling or enabling internally generated constraint triggers requires superuser privileges; it should be done with caution since of course the integrity of the constraint cannot be guaranteed if the triggers are not executed.

The trigger firing mechanism is also affected by the configuration variable [session_replication_role](#). Simply enabled triggers (the default) will fire when the replication role is “origin” (the default) or “local”. Triggers configured as `ENABLE REPLICA` will only fire if the session is in “replica” mode, and triggers configured as `ENABLE ALWAYS` will fire regardless of the current replication role.

The effect of this mechanism is that in the default configuration, triggers do not fire on replicas. This is useful because if a trigger is used on the origin to propagate data between tables, then the replication system will also replicate the propagated data; so the trigger should not fire a second time on the replica, because that would lead to duplication. However, if a trigger is used for another purpose such as creating external alerts, then it might be appropriate to set it to `ENABLE ALWAYS` so that it is also fired on replicas.

When this command is applied to a partitioned table, the states of corresponding clone triggers in the partitions are updated too, unless `ONLY` is specified.

This command acquires a `SHARE ROW EXCLUSIVE` lock.

DISABLE/ENABLE [REPLICA | ALWAYS] RULE

These forms configure the firing of rewrite rules belonging to the table. A disabled rule is still known to the system, but is not applied during query rewriting. The semantics are as for disabled/enabled triggers. This configuration is ignored for `ON SELECT` rules, which are always applied in order to keep views working even if the current session is in a non-default replication role.

The rule firing mechanism is also affected by the configuration variable [session_replication_role](#), analogous to triggers as described above.

DISABLE/ENABLE ROW LEVEL SECURITY

These forms control the application of row security policies belonging to the table. If enabled and no policies exist for the table, then a default-deny policy is applied. Note that policies can exist for a table even if row-level security is disabled. In this case, the policies will *not* be applied and the policies will be ignored. See also [CREATE POLICY](#).

NO FORCE/FORCE ROW LEVEL SECURITY

These forms control the application of row security policies belonging to the table when the user is the table owner. If enabled, row-level security policies will be applied when the user is the table owner. If disabled (the default) then row-level security will not be applied when the user is the table owner. See also [CREATE POLICY](#).

CLUSTER ON

This form selects the default index for future [CLUSTER](#) operations. It does not actually re-cluster the table.

Changing cluster options acquires a `SHARE UPDATE EXCLUSIVE` lock.

SET WITHOUT CLUSTER

This form removes the most recently used [CLUSTER](#) index specification from the table. This affects future cluster operations that don't specify an index.

Changing cluster options acquires a `SHARE UPDATE EXCLUSIVE` lock.

SET WITHOUT OIDS

Backward-compatible syntax for removing the `oid` system column. As `oid` system columns cannot be added anymore, this never has an effect.

SET ACCESS METHOD

This form changes the access method of the table by rewriting it using the indicated access method; specifying `DEFAULT` selects the access method set as the [default_table_access_method](#) configuration parameter. See [Chapter 62](#) for more information.

When applied to a partitioned table, there is no data to rewrite, but partitions created afterwards will default to the given access method unless overridden by a `USING` clause. Specifying `DEFAULT` removes a previous value, causing future partitions to default to `default_table_access_method`.

SET TABLESPACE

This form changes the table's tablespace to the specified tablespace and moves the data file(s) associated with the table to the new tablespace. Indexes on the table, if any, are not moved; but they can be moved separately with additional `SET TABLESPACE` commands. When applied to a partitioned table, nothing is moved, but any partitions created afterwards with `CREATE TABLE PARTITION OF` will use that tablespace, unless overridden by a `TABLESPACE` clause.

All tables in the current database in a tablespace can be moved by using the `ALL IN TABLESPACE` form, which will lock all tables to be moved first and then move each one. This form also supports `OWNED BY`, which will only move tables owned by the roles specified. If the `NOWAIT` option is specified then the command will fail if it is unable to acquire all of the locks required immediately. Note that system catalogs are not moved by this command; use `ALTER DATABASE` or explicit `ALTER TABLE` invocations instead if desired. The `information_schema` relations are not considered part of the system catalogs and will be moved. See also [CREATE TABLESPACE](#).

SET { LOGGED | UNLOGGED }

This form changes the table from unlogged to logged or vice-versa (see [UNLOGGED](#)). It cannot be applied to a temporary table.

This also changes the persistence of any sequences linked to the table (for identity or serial columns). However, it is also possible to change the persistence of such sequences separately.

SET CONSTANT

This form changes the table into read-only mode. No data can be modified or added to constant tables, and they are not processed by [autovacuum](#). Constant tables cannot be changed back to read-write mode.

SET (*storage_parameter* [= *value*] [, ...])

This form changes one or more storage parameters for the table. See [Storage Parameters](#) in the [CREATE TABLE](#) documentation for details on the available parameters. Note that the table contents will not be modified immediately by this command; depending on the parameter you might need to rewrite the table to get the desired effects. That can be done with [VACUUM FULL](#), [CLUSTER](#) or one of the forms of [ALTER TABLE](#) that forces a table rewrite. For planner related parameters, changes will take effect from the next time the table is locked so currently executing queries will not be affected.

SHARE UPDATE EXCLUSIVE lock will be taken for fillfactor, toast and autovacuum storage parameters, as well as the planner parameter `parallel_workers`.

RESET (*storage_parameter* [, ...])

This form resets one or more storage parameters to their defaults. As with `SET`, a table rewrite might be needed to update the table entirely.

INHERIT *parent_table*

This form adds the target table as a new child of the specified parent table. Subsequently, queries against the parent will include records of the target table. To be added as a child, the target table must already contain all the same columns as the parent (it could have additional columns, too). The columns must have matching data types, and if they have `NOT NULL` constraints in the parent then they must also have `NOT NULL` constraints in the child.

There must also be matching child-table constraints for all `CHECK` constraints of the parent, except those marked non-inheritable (that is, created with `ALTER TABLE ... ADD CONSTRAINT ... NO INHERIT`) in the parent, which are ignored; all child-table constraints matched must not be marked non-inheritable. Currently `UNIQUE`, `PRIMARY KEY`, and `FOREIGN KEY` constraints are not considered, but this might change in the future.

NO INHERIT *parent_table*

This form removes the target table from the list of children of the specified parent table. Queries against the parent table will no longer include records drawn from the target table.

OF *type_name*

This form links the table to a composite type as though `CREATE TABLE OF` had formed it. The table's list of column names and types must precisely match that of the composite type. The table must not inherit from any other table. These restrictions ensure that `CREATE TABLE OF` would permit an equivalent table definition.

NOT OF

This form dissociates a typed table from its type.

OWNER TO

This form changes the owner of the table, sequence, view, materialized view, or foreign table to the specified user.

REPLICA IDENTITY

This form changes the information which is written to the write-ahead log to identify rows which are updated or deleted. In most cases, the old value of each column is only logged if it differs from the

new value; however, if the old value is stored externally, it is always logged regardless of whether it changed. This option has no effect except when logical replication is in use.

DEFAULT

Records the old values of the columns of the primary key, if any. This is the default for non-system tables.

USING INDEX *index_name*

Records the old values of the columns covered by the named index, that must be unique, not partial, not deferrable, and include only columns marked `NOT NULL`. If this index is dropped, the behavior is the same as `NOTHING`.

FULL

Records the old values of all columns in the row.

NOTHING

Records no information about the old row. This is the default for system tables.

RENAME

The `RENAME` forms change the name of a table (or an index, sequence, view, materialized view, or foreign table), the name of an individual column in a table, or the name of a constraint of the table. When renaming a constraint that has an underlying index, the index is renamed as well. There is no effect on the stored data.

SET SCHEMA

This form moves the table into another schema. Associated indexes, constraints, and sequences owned by table columns are moved as well.

ATTACH PARTITION *partition_name* { FOR VALUES *partition_bound_spec* | DEFAULT }

This form attaches an existing table (which might itself be partitioned) as a partition of the target table. The table can be attached as a partition for specific values using `FOR VALUES` or as a default partition by using `DEFAULT`. For each index in the target table, a corresponding one will be created in the attached table; or, if an equivalent index already exists, it will be attached to the target table's index, as if `ALTER INDEX ATTACH PARTITION` had been executed. Note that if the existing table is a foreign table, it is currently not allowed to attach the table as a partition of the target table if there are `UNIQUE` indexes on the target table. For each user-defined row-level trigger that exists in the target table, a corresponding one is created in the attached table.

A partition using `FOR VALUES` uses same syntax for *partition_bound_spec* as `CREATE TABLE`. The partition bound specification must correspond to the partitioning strategy and partition key of the target table. The table to be attached must have all the same columns as the target table and no more; moreover, the column types must also match. Also, it must have all the `NOT NULL` and `CHECK` constraints of the target table, not marked `NO INHERIT`. Currently `FOREIGN KEY` constraints are not considered. `UNIQUE` and `PRIMARY KEY` constraints from the parent table will be created in the partition, if they don't already exist.

If the new partition is a regular table, a full table scan is performed to check that existing rows in the table do not violate the partition constraint. It is possible to avoid this scan by adding a valid `CHECK` constraint to the table that allows only rows satisfying the desired partition constraint before running this command. The `CHECK` constraint will be used to determine that the table need not be scanned to validate the partition constraint. This does not work, however, if any of the partition keys is an expression and the partition does not accept `NULL` values. If attaching a list partition that will not accept `NULL` values, also add a `NOT NULL` constraint to the partition key column, unless it's an expression.

If the new partition is a foreign table, nothing is done to verify that all the rows in the foreign table obey the partition constraint.

When a table has a default partition, defining a new partition changes the partition constraint for the default partition. The default partition can't contain any rows that would need to be moved to the new partition, and will be scanned to verify that none are present. This scan, like the scan of the new partition, can be avoided if an appropriate `CHECK` constraint is present. Also like the scan of the new partition, it is always skipped when the default partition is a foreign table.

Attaching a partition acquires a `SHARE UPDATE EXCLUSIVE` lock on the parent table, in addition to the `ACCESS EXCLUSIVE` locks on the table being attached and on the default partition (if any).

Further locks must also be held on all sub-partitions if the table being attached is itself a partitioned table. Likewise if the default partition is itself a partitioned table. The locking of the sub-partitions can be avoided by adding a `CHECK` constraint as described in [Section 5.12.2.2](#).

```
DETACH PARTITION partition_name [ CONCURRENTLY | FINALIZE ]
```

This form detaches the specified partition of the target table. The detached partition continues to exist as a standalone table, but no longer has any ties to the table from which it was detached. Any indexes that were attached to the target table's indexes are detached. Any triggers that were created as clones of those in the target table are removed. `SHARE` lock is obtained on any tables that reference this partitioned table in foreign key constraints.

If `CONCURRENTLY` is specified, it runs using a reduced lock level to avoid blocking other sessions that might be accessing the partitioned table. In this mode, two transactions are used internally. During the first transaction, a `SHARE UPDATE EXCLUSIVE` lock is taken on both parent table and partition, and the partition is marked as undergoing detach; at that point, the transaction is committed and all other transactions using the partitioned table are waited for. Once all those transactions have completed, the second transaction acquires `SHARE UPDATE EXCLUSIVE` on the partitioned table and `ACCESS EXCLUSIVE` on the partition, and the detach process completes. A `CHECK` constraint that duplicates the partition constraint is added to the partition. `CONCURRENTLY` cannot be run in a transaction block and is not allowed if the partitioned table contains a default partition.

If `FINALIZE` is specified, a previous `DETACH CONCURRENTLY` invocation that was canceled or interrupted is completed. At most one partition in a partitioned table can be pending detach at a time.

```
SPLIT PARTITION partition_name INTO (PARTITION partition_name1 { FOR VALUES partition_bound_spec | DEFAULT }, PARTITION partition_name2 { FOR VALUES partition_bound_spec | DEFAULT } [, ...])
```

This form splits a single partition of the target table. Hash-partitioning is not supported. Bounds of new partitions should not overlap with new and existing partitions (except *partition_name*). If the split partition is a `DEFAULT` partition, one of the new partitions must be `DEFAULT`. In case one of the new partitions or one of existing partitions is `DEFAULT`, new partitions *partition_name1*, *partition_name2*, ... can have spaces between partitions bounds. If the partitioned table does not have a `DEFAULT` partition, the `DEFAULT` partition can be defined as one of the new partitions.

In case new partitions do not contain a `DEFAULT` partition and the partitioned table does not have a `DEFAULT` partition, the following must be true: sum bounds of new partitions *partition_name1*, *partition_name2*, ... should be equal to bound of split partition *partition_name*. One of the new partitions *partition_name1*, *partition_name2*, ... can have the same name as split partition *partition_name* (this is suitable in case of splitting a `DEFAULT` partition: we split it, but after splitting we have a partition with the same name). Only simple, non-partitioned partition can be split.

The new partitions will be created the same as tables created with the SQL command `CREATE TABLE partition_nameN (LIKE name INCLUDING ALL EXCLUDING INDEXES EXCLUDING IDENTITY EXCLUDING STATISTICS)`. The indexes and identity are created later, after moving the data into the new partitions. Extended statistics aren't copied from the parent table, for consistency with `CREATE TABLE PARTITION OF`. New partitions will have the same table access method as the parent. If the parent table is persistent then new partitions are created persistent. If the parent table is temporary then new partitions are also created temporary. New partitions will also be created in the same tablespace as the parent.

Note

This command acquires an `ACCESS EXCLUSIVE` lock. This is a significant limitation, which limits the usage of this command with large partitioned tables under a high load.

```
MERGE PARTITIONS (partition_name1, partition_name2 [, ...]) INTO partition_name
```

This form merges several partitions into the one partition of the target table. Hash-partitioning is not supported. If `DEFAULT` partition is not in the list of partitions *partition_name1*, *partition_name2* [, ...]:

- For range-partitioned tables it is necessary that the ranges of the partitions *partition_name1*, *partition_name2* [, ...] can be merged into one range without spaces and overlaps (otherwise an error will be generated). The combined range will be the range for the partition *partition_name*.

- For list-partitioned tables the value lists of all partitions *partition_name1*, *partition_name2* [, ...] are combined and form the list of values of partition *partition_name*.

If `DEFAULT` partition is in the list of partitions *partition_name1*, *partition_name2* [, ...]:

- The partition *partition_name* will be the `DEFAULT` partition.
- For range- and list-partitioned tables the ranges and lists of values of the merged partitions can be any.

The new partition *partition_name* can have the same name as one of the merged partitions. Only simple, non-partitioned partitions can be merged.

The new partition will be created the same as a table created with the SQL command `CREATE TABLE partition_name (LIKE name INCLUDING ALL EXCLUDING INDEXES EXCLUDING IDENTITY EXCLUDING STATISTICS)`. The indexes and identity are created later, after moving the data into the new partition. Extended statistics aren't copied from the parent table, for consistency with `CREATE TABLE PARTITION OF`. The new partition will have the same table access method as the parent. If the parent table is persistent then the new partition is created persistent. If the parent table is temporary then the new partition is also created temporary. The new partition will also be created in the same tablespace as the parent.

Note

This command acquires an `ACCESS EXCLUSIVE` lock. This is a significant limitation, which limits the usage of this command with large partitioned tables under a high load.

All the forms of `ALTER TABLE` that act on a single table, except `RENAME`, `SET SCHEMA`, `ATTACH PARTITION`, `DETACH PARTITION`, `SPLIT PARTITION`, and `MERGE PARTITIONS`, can be combined into a list of multiple alterations to be applied together. For example, it is possible to add several columns and/or alter the type of several columns in a single command. This is particularly useful with large tables, since only one pass over the table need be made.

You must own the table to use `ALTER TABLE`. To change the schema or tablespace of a table, you must also have `CREATE` privilege on the new schema or tablespace. To add the table as a new child of a parent table, you must own the parent table as well. Also, to attach a table as a new partition of the table, you must own the table being attached. To alter the owner, you must be able to `SET ROLE` to the new owning role, and that role must have `CREATE` privilege on the table's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the table. However, a superuser can alter ownership of any table anyway.) To add a column or alter a column type or use the `OF` clause, you must also have `USAGE` privilege on the data type.

Parameters

IF EXISTS

Do not throw an error if the table does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing table to alter. If `ONLY` is specified before the table name, only that table is altered. If `ONLY` is not specified, the table and all its descendant tables (if any) are altered. Optionally, `*` can be specified after the table name to explicitly indicate that descendant tables are included.

column_name

Name of a new or existing column.

new_column_name

New name for an existing column.

new_name

New name for the table.

data_type

Data type of the new column, or new data type for an existing column.

table_constraint

New table constraint for the table.

constraint_name

Name of a new or existing constraint.

CASCADE

Automatically drop objects that depend on the dropped column or constraint (for example, views referencing the column), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the column or constraint if there are any dependent objects. This is the default behavior.

trigger_name

Name of a single trigger to disable or enable.

ALL

Disable or enable all triggers belonging to the table. (This requires superuser privilege if any of the triggers are internally generated constraint triggers, such as those that are used to implement foreign key constraints or deferrable uniqueness and exclusion constraints.)

USER

Disable or enable all triggers belonging to the table except for internally generated constraint triggers, such as those that are used to implement foreign key constraints or deferrable uniqueness and exclusion constraints.

index_name

The name of an existing index.

storage_parameter

The name of a table storage parameter.

value

The new value for a table storage parameter. This might be a number or a word depending on the parameter.

parent_table

A parent table to associate or de-associate with this table.

new_owner

The user name of the new owner of the table.

new_access_method

The name of the access method to which the table will be converted.

new_tablespace

The name of the tablespace to which the table will be moved.

new_schema

The name of the schema to which the table will be moved.

partition_name

The name of the table to attach as a new partition or to detach from this table, or the name of split partition, or the name of the new merged partition.

partition_bound_spec

The partition bound specification for a new partition. Refer to [CREATE TABLE](#) for more details on the syntax of the same.

Storage Parameters

Postgres Pro Shardman extends *storage parameters* of tables with its own *storage metaparameters*. They are not stored in the corresponding catalog entry, but are used to tell the Postgres Pro Shardman extension to perform some additional actions.

global

This parameter can be specified only for global tables. If set to 0, the global table will be converted to a regular one on the replication group where the command is executed. The global table will not exist on other nodes after completion of this statement. However, setting storage options is allowed for global tables.

Notes

The key word `COLUMN` is noise and can be omitted.

When a column is added with `ADD COLUMN` and a non-volatile `DEFAULT` is specified, the default is evaluated at the time of the statement and the result stored in the table's metadata. That value will be used for the column for all existing rows. If no `DEFAULT` is specified, `NULL` is used. In neither case is a rewrite of the table required.

Adding a column with a volatile `DEFAULT` or changing the type of an existing column will require the entire table and its indexes to be rewritten. As an exception, when changing the type of an existing column, if the `USING` clause does not change the column contents and the old type is either binary coercible to the new type or an unconstrained domain over the new type, a table rewrite is not needed. However, indexes

must always be rebuilt unless the system can verify that the new index would be logically equivalent to the existing one. For example, if the collation for a column has been changed, an index rebuild is always required because the new sort order might be different. However, in the absence of a collation change, a column can be changed from `text` to `varchar` (or vice versa) without rebuilding the indexes because these data types sort identically. Table and/or index rebuilds may take a significant amount of time for a large table; and will temporarily require as much as double the disk space.

Adding a `CHECK` or `NOT NULL` constraint requires scanning the table to verify that existing rows meet the constraint, but does not require a table rewrite.

Similarly, when attaching a new partition it may be scanned to verify that existing rows meet the partition constraint.

The main reason for providing the option to specify multiple changes in a single `ALTER TABLE` is that multiple table scans or rewrites can thereby be combined into a single pass over the table.

Scanning a large table to verify a new foreign key or check constraint can take a long time, and other updates to the table are locked out until the `ALTER TABLE ADD CONSTRAINT` command is committed. The main purpose of the `NOT VALID` constraint option is to reduce the impact of adding a constraint on concurrent updates. With `NOT VALID`, the `ADD CONSTRAINT` command does not scan the table and can be committed immediately. After that, a `VALIDATE CONSTRAINT` command can be issued to verify that existing rows satisfy the constraint. The validation step does not need to lock out concurrent updates, since it knows that other transactions will be enforcing the constraint for rows that they insert or update; only pre-existing rows need to be checked. Hence, validation acquires only a `SHARE UPDATE EXCLUSIVE` lock on the table being altered. (If the constraint is a foreign key then a `ROW SHARE` lock is also required on the table referenced by the constraint.) In addition to improving concurrency, it can be useful to use `NOT VALID` and `VALIDATE CONSTRAINT` in cases where the table is known to contain pre-existing violations. Once the constraint is in place, no new violations can be inserted, and the existing problems can be corrected at leisure until `VALIDATE CONSTRAINT` finally succeeds.

The `DROP COLUMN` form does not physically remove the column, but simply makes it invisible to SQL operations. Subsequent insert and update operations in the table will store a null value for the column. Thus, dropping a column is quick but it will not immediately reduce the on-disk size of your table, as the space occupied by the dropped column is not reclaimed. The space will be reclaimed over time as existing rows are updated.

To force immediate reclamation of space occupied by a dropped column, you can execute one of the forms of `ALTER TABLE` that performs a rewrite of the whole table. This results in reconstructing each row with the dropped column replaced by a null value.

The rewriting forms of `ALTER TABLE` are not MVCC-safe. After a table rewrite, the table will appear empty to concurrent transactions, if they are using a snapshot taken before the rewrite occurred. See [Section 13.6](#) for more details.

The `USING` option of `SET DATA TYPE` can actually specify any expression involving the old values of the row; that is, it can refer to other columns as well as the one being converted. This allows very general conversions to be done with the `SET DATA TYPE` syntax. Because of this flexibility, the `USING` expression is not applied to the column's default value (if any); the result might not be a constant expression as required for a default. This means that when there is no implicit or assignment cast from old to new type, `SET DATA TYPE` might fail to convert the default even though a `USING` clause is supplied. In such cases, drop the default with `DROP DEFAULT`, perform the `ALTER TYPE`, and then use `SET DEFAULT` to add a suitable new default. Similar considerations apply to indexes and constraints involving the column.

If a table has any descendant tables, it is not permitted to add, rename, or change the type of a column in the parent table without doing the same to the descendants. This ensures that the descendants always have columns matching the parent. Similarly, a `CHECK` constraint cannot be renamed in the parent without also renaming it in all descendants, so that `CHECK` constraints also match between the parent and its descendants. (That restriction does not apply to index-based constraints, however.) Also, because selecting from the parent also selects from its descendants, a constraint on the parent cannot be marked

valid unless it is also marked valid for those descendants. In all of these cases, `ALTER TABLE ONLY` will be rejected.

A recursive `DROP COLUMN` operation will remove a descendant table's column only if the descendant does not inherit that column from any other parents and never had an independent definition of the column. A nonrecursive `DROP COLUMN` (i.e., `ALTER TABLE ONLY ... DROP COLUMN`) never removes any descendant columns, but instead marks them as independently defined rather than inherited. A nonrecursive `DROP COLUMN` command will fail for a partitioned table, because all partitions of a table must have the same columns as the partitioning root.

The actions for identity columns (`ADD GENERATED`, `SET` etc., `DROP IDENTITY`), as well as the actions `CLUSTER`, `OWNER`, and `TABLESPACE` never recurse to descendant tables; that is, they always act as though `ONLY` were specified. Actions affecting trigger states recurse to partitions of partitioned tables (unless `ONLY` is specified), but never to traditional-inheritance descendants. Adding a constraint recurses only for `CHECK` constraints that are not marked `NO INHERIT`.

Changing any part of a system catalog table is not permitted.

Refer to [CREATE TABLE](#) for a further description of valid parameters. [Chapter 5](#) has further information on inheritance.

Examples

To add a column of type `varchar` to a table:

```
ALTER TABLE distributors ADD COLUMN address varchar(30);
```

That will cause all existing rows in the table to be filled with null values for the new column.

To add a column with a non-null default:

```
ALTER TABLE measurements
  ADD COLUMN mtime timestamp with time zone DEFAULT now();
```

Existing rows will be filled with the current time as the value of the new column, and then new rows will receive the time of their insertion.

To add a column and fill it with a value different from the default to be used later:

```
ALTER TABLE transactions
  ADD COLUMN status varchar(30) DEFAULT 'old',
  ALTER COLUMN status SET default 'current';
```

Existing rows will be filled with `old`, but then the default for subsequent commands will be `current`. The effects are the same as if the two sub-commands had been issued in separate `ALTER TABLE` commands.

To drop a column from a table:

```
ALTER TABLE distributors DROP COLUMN address RESTRICT;
```

To change the types of two existing columns in one operation:

```
ALTER TABLE distributors
  ALTER COLUMN address TYPE varchar(80),
  ALTER COLUMN name TYPE varchar(100);
```

To change an integer column containing Unix timestamps to `timestamp with time zone` via a `USING` clause:

```
ALTER TABLE foo
  ALTER COLUMN foo_timestamp SET DATA TYPE timestamp with time zone
  USING
    timestamp with time zone 'epoch' + foo_timestamp * interval '1 second';
```

The same, when the column has a default expression that won't automatically cast to the new data type:


```
ALTER TABLE foo
  ALTER COLUMN foo_timestamp DROP DEFAULT,
  ALTER COLUMN foo_timestamp TYPE timestamp with time zone
  USING
    timestamp with time zone 'epoch' + foo_timestamp * interval '1 second',
  ALTER COLUMN foo_timestamp SET DEFAULT now();
```

To rename an existing column:

```
ALTER TABLE distributors RENAME COLUMN address TO city;
```

To rename an existing table:

```
ALTER TABLE distributors RENAME TO suppliers;
```

To rename an existing constraint:

```
ALTER TABLE distributors RENAME CONSTRAINT zipchk TO zip_check;
```

To add a not-null constraint to a column:

```
ALTER TABLE distributors ALTER COLUMN street SET NOT NULL;
```

To remove a not-null constraint from a column:

```
ALTER TABLE distributors ALTER COLUMN street DROP NOT NULL;
```

To add a check constraint to a table and all its children:

```
ALTER TABLE distributors ADD CONSTRAINT zipchk CHECK (char_length(zipcode) = 5);
```

To add a check constraint only to a table and not to its children:

```
ALTER TABLE distributors ADD CONSTRAINT zipchk CHECK (char_length(zipcode) = 5) NO
  INHERIT;
```

(The check constraint will not be inherited by future children, either.)

To remove a check constraint from a table and all its children:

```
ALTER TABLE distributors DROP CONSTRAINT zipchk;
```

To remove a check constraint from one table only:

```
ALTER TABLE ONLY distributors DROP CONSTRAINT zipchk;
```

(The check constraint remains in place for any child tables.)

To add a foreign key constraint to a table:

```
ALTER TABLE distributors ADD CONSTRAINT distfk FOREIGN KEY (address) REFERENCES
  addresses (address);
```

To add a foreign key constraint to a table with the least impact on other work:

```
ALTER TABLE distributors ADD CONSTRAINT distfk FOREIGN KEY (address) REFERENCES
  addresses (address) NOT VALID;
ALTER TABLE distributors VALIDATE CONSTRAINT distfk;
```

To add a (multicolumn) unique constraint to a table:

```
ALTER TABLE distributors ADD CONSTRAINT dist_id_zipcode_key UNIQUE (dist_id, zipcode);
```

To add an automatically named primary key constraint to a table, noting that a table can only ever have one primary key:

```
ALTER TABLE distributors ADD PRIMARY KEY (dist_id);
```

To move a table to a different tablespace:

```
ALTER TABLE distributors SET TABLESPACE fasttablespace;
```

To move a table to a different schema:

```
ALTER TABLE myschema.distributors SET SCHEMA yourschema;
```

To recreate a primary key constraint, without blocking updates while the index is rebuilt:

```
CREATE UNIQUE INDEX CONCURRENTLY dist_id_temp_idx ON distributors (dist_id);
ALTER TABLE distributors DROP CONSTRAINT distributors_pkey,
    ADD CONSTRAINT distributors_pkey PRIMARY KEY USING INDEX dist_id_temp_idx;
```

To attach a partition to a range-partitioned table:

```
ALTER TABLE measurement
    ATTACH PARTITION measurement_y2016m07 FOR VALUES FROM ('2016-07-01') TO
    ('2016-08-01');
```

To attach a partition to a list-partitioned table:

```
ALTER TABLE cities
    ATTACH PARTITION cities_ab FOR VALUES IN ('a', 'b');
```

To attach a partition to a hash-partitioned table:

```
ALTER TABLE orders
    ATTACH PARTITION orders_p4 FOR VALUES WITH (MODULUS 4, REMAINDER 3);
```

To attach a default partition to a partitioned table:

```
ALTER TABLE cities
    ATTACH PARTITION cities_partdef DEFAULT;
```

To detach a partition from a partitioned table:

```
ALTER TABLE measurement
    DETACH PARTITION measurement_y2015m12;
```

To split a single partition of the range-partitioned table:

```
ALTER TABLE sales_range SPLIT PARTITION sales_feb_mar_apr2023 INTO
    (PARTITION sales_feb2023 FOR VALUES FROM ('2023-02-01') TO ('2023-03-01'),
    PARTITION sales_mar2023 FOR VALUES FROM ('2023-03-01') TO ('2023-04-01'),
    PARTITION sales_apr2023 FOR VALUES FROM ('2023-04-01') TO ('2023-05-01'));
```

To split a single partition of the list-partitioned table:

```
ALTER TABLE sales_list SPLIT PARTITION sales_all INTO
    (PARTITION sales_west FOR VALUES IN ('Lisbon', 'New York', 'Madrid'),
    PARTITION sales_east FOR VALUES IN ('Beijing', 'Delhi', 'Vladivostok'),
    PARTITION sales_central FOR VALUES IN ('Warsaw', 'Berlin', 'Kyiv'));
```

To merge several partitions into one partition of the target table:

```
ALTER TABLE sales_list MERGE PARTITIONS (sales_west, sales_east, sales_central)
    INTO sales_all;
```

To alter a global table, create a global table pgbench_tellers and then convert it to local.

```
CREATE TABLE pgbench_tellers (
    tid        integer PRIMARY KEY,
    bid        integer,
    tbalance   integer,
    filler     character(84)
)
```

```
WITH (global);  
ALTER TABLE pgbench_tellers SET (global=0);
```

Compatibility

The forms `ADD (without USING INDEX)`, `DROP [COLUMN]`, `DROP IDENTITY`, `RESTART`, `SET DEFAULT`, `SET DATA TYPE (without USING)`, `SET GENERATED`, and `SET sequence_option` conform with the SQL standard. The other forms are Postgres Pro extensions of the SQL standard. Also, the ability to specify more than one manipulation in a single `ALTER TABLE` command is an extension.

`ALTER TABLE DROP COLUMN` can be used to drop the only column of a table, leaving a zero-column table. This is an extension of SQL, which disallows zero-column tables.

ALTER TABLE Limitations for Postgres Pro Shardman

- `ALTER TABLE` is prohibited for partitions of sharded tables.
- `ALTER TABLE OWNER` is allowed. For sharded table it also changes the owner of table partitions. Only the global user can be an owner of sharded or global table.
- `ALTER TABLE COLUMN TYPE` is allowed with limitations. You cannot alter type of sharded table column participating in sharding or partitioning key. You cannot alter type of sharded table column with `USING` clause (but for global tables it is allowed). Also, it is a user's duty for now to create and keep new type exactly equal on every cluster node.
- Adding or dropping table-wide unique constraints and checks is allowed. For global tables dropping primary key constraint or dropping columns, participating in primary key, is forbidden.
- Adding foreign keys between sharded tables is possible only when they are colocated and a foreign key references tuples that are stored in the same replication group. A foreign key between sharded tables must begin with the columns used for table partitioning in both tables. A foreign key on a global table can reference only another global table. A foreign key on a sharded table can reference a colocated sharded table or a global table.
- `SET/DROP NOT NULL` is allowed.
- If the parent table is global, the `ON DELETE` can only perform `NO ACTION`, `RESTRICT`, `CASCADE`, and `ON UPDATE` can only perform `NO ACTION`, `RESTRICT`.
- Setting storage options is allowed for global tables.
- Global and sharded tables cannot inherit other tables.
- `ALTER COLUMN SET STATISTICS` is allowed for global and sharded tables.

See Also

[CREATE TABLE](#)

ALTER TABLESPACE

ALTER TABLESPACE — change the definition of a tablespace

Synopsis

```
ALTER TABLESPACE name RENAME TO new_name
ALTER TABLESPACE name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER TABLESPACE name SET ( tablespace_option = value [, ... ] )
ALTER TABLESPACE name RESET ( tablespace_option [, ... ] )
```

Description

ALTER TABLESPACE can be used to change the definition of a tablespace.

You must own the tablespace to change the definition of a tablespace. To alter the owner, you must also be able to SET ROLE to the new owning role. (Note that superusers have these privileges automatically.)

Parameters

name

The name of an existing tablespace.

new_name

The new name of the tablespace. The new name cannot begin with `pg_`, as such names are reserved for system tablespaces.

new_owner

The new owner of the tablespace.

tablespace_option

A tablespace parameter to be set or reset. Currently, the only available parameters are `seq_page_cost`, `random_page_cost`, `effective_io_concurrency` and `maintenance_io_concurrency`. Setting these values for a particular tablespace will override the planner's usual estimate of the cost of reading pages from tables in that tablespace, and the executor's prefetching behavior, as established by the configuration parameters of the same name (see [seq_page_cost](#), [random_page_cost](#), [effective_io_concurrency](#), [maintenance_io_concurrency](#)). This may be useful if one tablespace is located on a disk which is faster or slower than the remainder of the I/O subsystem.

Examples

Rename tablespace `index_space` to `fast_raid`:

```
ALTER TABLESPACE index_space RENAME TO fast_raid;
```

Change the owner of tablespace `index_space`:

```
ALTER TABLESPACE index_space OWNER TO mary;
```

Compatibility

There is no ALTER TABLESPACE statement in the SQL standard.

See Also

[CREATE TABLESPACE](#), [DROP TABLESPACE](#)

ALTER TEXT SEARCH CONFIGURATION

ALTER TEXT SEARCH CONFIGURATION — change the definition of a text search configuration

Synopsis

```
ALTER TEXT SEARCH CONFIGURATION name
    ADD MAPPING FOR token_type [, ... ] WITH dictionary_name [, ... ]
ALTER TEXT SEARCH CONFIGURATION name
    ALTER MAPPING FOR token_type [, ... ] WITH dictionary_name [, ... ]
ALTER TEXT SEARCH CONFIGURATION name
    ALTER MAPPING REPLACE old_dictionary WITH new_dictionary
ALTER TEXT SEARCH CONFIGURATION name
    ALTER MAPPING FOR token_type [, ... ] REPLACE old_dictionary WITH new_dictionary
ALTER TEXT SEARCH CONFIGURATION name
    DROP MAPPING [ IF EXISTS ] FOR token_type [, ... ]
ALTER TEXT SEARCH CONFIGURATION name RENAME TO new_name
ALTER TEXT SEARCH CONFIGURATION name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER
    | SESSION_USER }
ALTER TEXT SEARCH CONFIGURATION name SET SCHEMA new_schema
```

Description

ALTER TEXT SEARCH CONFIGURATION changes the definition of a text search configuration. You can modify its mappings from token types to dictionaries, or change the configuration's name or owner.

You must be the owner of the configuration to use ALTER TEXT SEARCH CONFIGURATION.

Parameters

name

The name (optionally schema-qualified) of an existing text search configuration.

token_type

The name of a token type that is emitted by the configuration's parser.

dictionary_name

The name of a text search dictionary to be consulted for the specified token type(s). If multiple dictionaries are listed, they are consulted in the specified order.

old_dictionary

The name of a text search dictionary to be replaced in the mapping.

new_dictionary

The name of a text search dictionary to be substituted for *old_dictionary*.

new_name

The new name of the text search configuration.

new_owner

The new owner of the text search configuration.

new_schema

The new schema for the text search configuration.

The `ADD MAPPING FOR` form installs a list of dictionaries to be consulted for the specified token type(s); it is an error if there is already a mapping for any of the token types. The `ALTER MAPPING FOR` form does the same, but first removing any existing mapping for those token types. The `ALTER MAPPING REPLACE` forms substitute *new_dictionary* for *old_dictionary* anywhere the latter appears. This is done for only the specified token types when `FOR` appears, or for all mappings of the configuration when it doesn't. The `DROP MAPPING` form removes all dictionaries for the specified token type(s), causing tokens of those types to be ignored by the text search configuration. It is an error if there is no mapping for the token types, unless `IF EXISTS` appears.

Examples

The following example replaces the `english` dictionary with the `swedish` dictionary anywhere that `english` is used within `my_config`.

```
ALTER TEXT SEARCH CONFIGURATION my_config
    ALTER MAPPING REPLACE english WITH swedish;
```

Compatibility

There is no `ALTER TEXT SEARCH CONFIGURATION` statement in the SQL standard.

See Also

[CREATE TEXT SEARCH CONFIGURATION](#), [DROP TEXT SEARCH CONFIGURATION](#)

ALTER TEXT SEARCH DICTIONARY

ALTER TEXT SEARCH DICTIONARY — change the definition of a text search dictionary

Synopsis

```
ALTER TEXT SEARCH DICTIONARY name (  
    option [ = value ] [, ... ]  
)  
ALTER TEXT SEARCH DICTIONARY name RENAME TO new_name  
ALTER TEXT SEARCH DICTIONARY name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |  
    SESSION_USER }  
ALTER TEXT SEARCH DICTIONARY name SET SCHEMA new_schema
```

Description

ALTER TEXT SEARCH DICTIONARY changes the definition of a text search dictionary. You can change the dictionary's template-specific options, or change the dictionary's name or owner.

You must be the owner of the dictionary to use ALTER TEXT SEARCH DICTIONARY.

Parameters

name

The name (optionally schema-qualified) of an existing text search dictionary.

option

The name of a template-specific option to be set for this dictionary.

value

The new value to use for a template-specific option. If the equal sign and value are omitted, then any previous setting for the option is removed from the dictionary, allowing the default to be used.

new_name

The new name of the text search dictionary.

new_owner

The new owner of the text search dictionary.

new_schema

The new schema for the text search dictionary.

Template-specific options can appear in any order.

Examples

The following example command changes the stopword list for a Snowball-based dictionary. Other parameters remain unchanged.

```
ALTER TEXT SEARCH DICTIONARY my_dict ( StopWords = newrussian );
```

The following example command changes the language option to `dutch`, and removes the stopword option entirely.

```
ALTER TEXT SEARCH DICTIONARY my_dict ( language = dutch, StopWords );
```

The following example command “updates” the dictionary's definition without actually changing anything.

```
ALTER TEXT SEARCH DICTIONARY my_dict ( dummy );
```

(The reason this works is that the option removal code doesn't complain if there is no such option.) This trick is useful when changing configuration files for the dictionary: the `ALTER` will force existing database sessions to re-read the configuration files, which otherwise they would never do if they had read them earlier.

Compatibility

There is no `ALTER TEXT SEARCH DICTIONARY` statement in the SQL standard.

See Also

[CREATE TEXT SEARCH DICTIONARY](#), [DROP TEXT SEARCH DICTIONARY](#)

ALTER TEXT SEARCH PARSER

ALTER TEXT SEARCH PARSER — change the definition of a text search parser

Synopsis

```
ALTER TEXT SEARCH PARSER name RENAME TO new_name
ALTER TEXT SEARCH PARSER name SET SCHEMA new_schema
```

Description

ALTER TEXT SEARCH PARSER changes the definition of a text search parser. Currently, the only supported functionality is to change the parser's name.

You must be a superuser to use ALTER TEXT SEARCH PARSER.

Parameters

name

The name (optionally schema-qualified) of an existing text search parser.

new_name

The new name of the text search parser.

new_schema

The new schema for the text search parser.

Compatibility

There is no ALTER TEXT SEARCH PARSER statement in the SQL standard.

See Also

[CREATE TEXT SEARCH PARSER](#), [DROP TEXT SEARCH PARSER](#)

ALTER TEXT SEARCH TEMPLATE

ALTER TEXT SEARCH TEMPLATE — change the definition of a text search template

Synopsis

```
ALTER TEXT SEARCH TEMPLATE name RENAME TO new_name
ALTER TEXT SEARCH TEMPLATE name SET SCHEMA new_schema
```

Description

ALTER TEXT SEARCH TEMPLATE changes the definition of a text search template. Currently, the only supported functionality is to change the template's name.

You must be a superuser to use ALTER TEXT SEARCH TEMPLATE.

Parameters

name

The name (optionally schema-qualified) of an existing text search template.

new_name

The new name of the text search template.

new_schema

The new schema for the text search template.

Compatibility

There is no ALTER TEXT SEARCH TEMPLATE statement in the SQL standard.

See Also

[CREATE TEXT SEARCH TEMPLATE](#), [DROP TEXT SEARCH TEMPLATE](#)

ALTER TRIGGER

ALTER TRIGGER — change the definition of a trigger

Synopsis

```
ALTER TRIGGER name ON table_name RENAME TO new_name
ALTER TRIGGER name ON table_name [ NO ] DEPENDS ON EXTENSION extension_name
```

Description

ALTER TRIGGER changes properties of an existing trigger.

The RENAME clause changes the name of the given trigger without otherwise changing the trigger definition. If the table that the trigger is on is a partitioned table, then corresponding clone triggers in the partitions are renamed too.

The DEPENDS ON EXTENSION clause marks the trigger as dependent on an extension, such that if the extension is dropped, the trigger will automatically be dropped as well.

You must own the table on which the trigger acts to be allowed to change its properties.

Parameters

name

The name of an existing trigger to alter.

table_name

The name of the table on which this trigger acts.

new_name

The new name for the trigger.

extension_name

The name of the extension that the trigger is to depend on (or no longer dependent on, if NO is specified). A trigger that's marked as dependent on an extension is automatically dropped when the extension is dropped.

Notes

The ability to temporarily enable or disable a trigger is provided by [ALTER TABLE](#), not by ALTER TRIGGER, because ALTER TRIGGER has no convenient way to express the option of enabling or disabling all of a table's triggers at once.

Examples

To rename an existing trigger:

```
ALTER TRIGGER emp_stamp ON emp RENAME TO emp_track_chgs;
```

To mark a trigger as being dependent on an extension:

```
ALTER TRIGGER emp_stamp ON emp DEPENDS ON EXTENSION emplib;
```

Compatibility

ALTER TRIGGER is a Postgres Pro extension of the SQL standard.

See Also

[ALTER TABLE](#)

ALTER TYPE

ALTER TYPE — change the definition of a type

Synopsis

```
ALTER TYPE name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
ALTER TYPE name RENAME TO new_name
ALTER TYPE name SET SCHEMA new_schema
ALTER TYPE name RENAME ATTRIBUTE attribute_name TO new_attribute_name [ CASCADE |
    RESTRICT ]
ALTER TYPE name action [, ... ]
ALTER TYPE name ADD VALUE [ IF NOT EXISTS ] new_enum_value [ { BEFORE |
    AFTER } neighbor_enum_value ]
ALTER TYPE name RENAME VALUE existing_enum_value TO new_enum_value
ALTER TYPE name SET ( property = value [, ... ] )
```

where *action* is one of:

```
    ADD ATTRIBUTE attribute_name data_type [ COLLATE collation ] [ CASCADE | RESTRICT ]
    DROP ATTRIBUTE [ IF EXISTS ] attribute_name [ CASCADE | RESTRICT ]
    ALTER ATTRIBUTE attribute_name [ SET DATA ] TYPE data_type [ COLLATE collation ]
[ CASCADE | RESTRICT ]
```

Description

ALTER TYPE changes the definition of an existing type. There are several subforms:

OWNER

This form changes the owner of the type.

RENAME

This form changes the name of the type.

SET SCHEMA

This form moves the type into another schema.

RENAME ATTRIBUTE

This form is only usable with composite types. It changes the name of an individual attribute of the type.

ADD ATTRIBUTE

This form adds a new attribute to a composite type, using the same syntax as [CREATE TYPE](#).

DROP ATTRIBUTE [IF EXISTS]

This form drops an attribute from a composite type. If IF EXISTS is specified and the attribute does not exist, no error is thrown. In this case a notice is issued instead.

ALTER ATTRIBUTE ... SET DATA TYPE

This form changes the type of an attribute of a composite type.

ADD VALUE [IF NOT EXISTS] [BEFORE | AFTER]

This form adds a new value to an enum type. The new value's place in the enum's ordering can be specified as being BEFORE or AFTER one of the existing values. Otherwise, the new item is added at the end of the list of values.

If `IF NOT EXISTS` is specified, it is not an error if the type already contains the new value: a notice is issued but no other action is taken. Otherwise, an error will occur if the new value is already present.

RENAME VALUE

This form renames a value of an enum type. The value's place in the enum's ordering is not affected. An error will occur if the specified value is not present or the new name is already present.

`SET (property = value [, ...])`

This form is only applicable to base types. It allows adjustment of a subset of the base-type properties that can be set in `CREATE TYPE`. Specifically, these properties can be changed:

- `RECEIVE` can be set to the name of a binary input function, or `NONE` to remove the type's binary input function. Using this option requires superuser privilege.
- `SEND` can be set to the name of a binary output function, or `NONE` to remove the type's binary output function. Using this option requires superuser privilege.
- `TYPMOD_IN` can be set to the name of a type modifier input function, or `NONE` to remove the type's type modifier input function. Using this option requires superuser privilege.
- `TYPMOD_OUT` can be set to the name of a type modifier output function, or `NONE` to remove the type's type modifier output function. Using this option requires superuser privilege.
- `ANALYZE` can be set to the name of a type-specific statistics collection function, or `NONE` to remove the type's statistics collection function. Using this option requires superuser privilege.
- `SUBSCRIPT` can be set to the name of a type-specific subscripting handler function, or `NONE` to remove the type's subscripting handler function. Using this option requires superuser privilege.
- `STORAGE` can be set to `plain`, `extended`, `external`, or `main` (see [Section 66.2](#) for more information about what these mean). However, changing from `plain` to another setting requires superuser privilege (because it requires that the type's C functions all be TOAST-ready), and changing to `plain` from another setting is not allowed at all (since the type may already have TOAST-ed values present in the database). Note that changing this option doesn't by itself change any stored data, it just sets the default TOAST strategy to be used for table columns created in the future. See [ALTER TABLE](#) to change the TOAST strategy for existing table columns.

See [CREATE TYPE](#) for more details about these type properties. Note that where appropriate, a change in these properties for a base type will be propagated automatically to domains based on that type.

The `ADD ATTRIBUTE`, `DROP ATTRIBUTE`, and `ALTER ATTRIBUTE` actions can be combined into a list of multiple alterations to apply in parallel. For example, it is possible to add several attributes and/or alter the type of several attributes in a single command.

You must own the type to use `ALTER TYPE`. To change the schema of a type, you must also have `CREATE` privilege on the new schema. To alter the owner, you must be able to `SET ROLE` to the new owning role, and that role must have `CREATE` privilege on the type's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the type. However, a superuser can alter ownership of any type anyway.) To add an attribute or alter an attribute type, you must also have `USAGE` privilege on the attribute's data type.

Parameters

name

The name (possibly schema-qualified) of an existing type to alter.

new_name

The new name for the type.

new_owner

The user name of the new owner of the type.

new_schema

The new schema for the type.

attribute_name

The name of the attribute to add, alter, or drop.

new_attribute_name

The new name of the attribute to be renamed.

data_type

The data type of the attribute to add, or the new type of the attribute to alter.

new_enum_value

The new value to be added to an enum type's list of values, or the new name to be given to an existing value. Like all enum literals, it needs to be quoted.

neighbor_enum_value

The existing enum value that the new value should be added immediately before or after in the enum type's sort ordering. Like all enum literals, it needs to be quoted.

existing_enum_value

The existing enum value that should be renamed. Like all enum literals, it needs to be quoted.

property

The name of a base-type property to be modified; see above for possible values.

CASCADE

Automatically propagate the operation to typed tables of the type being altered, and their descendants.

RESTRICT

Refuse the operation if the type being altered is the type of a typed table. This is the default.

Notes

If `ALTER TYPE ... ADD VALUE` (the form that adds a new value to an enum type) is executed inside a transaction block, the new value cannot be used until after the transaction has been committed.

Comparisons involving an added enum value will sometimes be slower than comparisons involving only original members of the enum type. This will usually only occur if `BEFORE` or `AFTER` is used to set the new value's sort position somewhere other than at the end of the list. However, sometimes it will happen even though the new value is added at the end (this occurs if the OID counter “wrapped around” since the original creation of the enum type). The slowdown is usually insignificant; but if it matters, optimal performance can be regained by dropping and recreating the enum type, or by dumping and restoring the database.

Examples

To rename a data type:

```
ALTER TYPE electronic_mail RENAME TO email;
```

To change the owner of the type `email` to `joe`:

```
ALTER TYPE email OWNER TO joe;
```

To change the schema of the type email to customers:

```
ALTER TYPE email SET SCHEMA customers;
```

To add a new attribute to a composite type:

```
ALTER TYPE compfoo ADD ATTRIBUTE f3 int;
```

To add a new value to an enum type in a particular sort position:

```
ALTER TYPE colors ADD VALUE 'orange' AFTER 'red';
```

To rename an enum value:

```
ALTER TYPE colors RENAME VALUE 'purple' TO 'mauve';
```

To create binary I/O functions for an existing base type:

```
CREATE FUNCTION mytypesend(mytype) RETURNS bytea ...;
CREATE FUNCTION mytyperecv(internal, oid, integer) RETURNS mytype ...;
ALTER TYPE mytype SET (
    SEND = mytypesend,
    RECEIVE = mytyperecv
);
```

Compatibility

The variants to add and drop attributes are part of the SQL standard; the other variants are Postgres Pro extensions.

See Also

[CREATE TYPE](#), [DROP TYPE](#)

ALTER USER

ALTER USER — change a database role

Synopsis

```
ALTER USER role_specification [ WITH ] option [ ... ]
```

where *option* can be:

```
SUPERUSER | NOSUPERUSER
| CREATEDB | NOCREATEDB
| CREATEROLE | NOCREATEROLE
| INHERIT | NOINHERIT
| LOGIN | NOLOGIN
| REPLICATION | NOREPLICATION
| BYPASSRLS | NOBYPASSRLS
| CONNECTION LIMIT connlimit
| [ ENCRYPTED ] PASSWORD 'password' | PASSWORD NULL
| VALID UNTIL 'timestamp'
| PROFILE profile_name
| ACCOUNT UNLOCK | ACCOUNT LOCK
```

```
ALTER USER name RENAME TO new_name
```

```
ALTER USER { role_specification | ALL } [ IN DATABASE database_name ]
SET configuration_parameter { TO | = } { value | DEFAULT }
ALTER USER { role_specification | ALL } [ IN DATABASE database_name ]
SET configuration_parameter FROM CURRENT
ALTER USER { role_specification | ALL } [ IN DATABASE database_name ]
RESET configuration_parameter
ALTER USER { role_specification | ALL } [ IN DATABASE database_name ] RESET ALL
```

where *role_specification* can be:

```
role_name
| CURRENT_ROLE
| CURRENT_USER
| SESSION_USER
```

Description

ALTER USER is now an alias for [ALTER ROLE](#).

Compatibility

The ALTER USER statement is a Postgres Pro extension. The SQL standard leaves the definition of users to the implementation.

See Also

[ALTER ROLE](#)

ALTER VIEW

ALTER VIEW — change the definition of a view

Synopsis

```
ALTER VIEW [ IF EXISTS ] name ALTER [ COLUMN ] column_name SET DEFAULT expression
ALTER VIEW [ IF EXISTS ] name ALTER [ COLUMN ] column_name DROP DEFAULT
ALTER VIEW [ IF EXISTS ] name OWNER TO { new_owner | CURRENT_ROLE | CURRENT_USER |
SESSION_USER }
ALTER VIEW [ IF EXISTS ] name RENAME [ COLUMN ] column_name TO new_column_name
ALTER VIEW [ IF EXISTS ] name RENAME TO new_name
ALTER VIEW [ IF EXISTS ] name SET SCHEMA new_schema
ALTER VIEW [ IF EXISTS ] name SET ( view_option_name [= view_option_value] [, ... ] )
ALTER VIEW [ IF EXISTS ] name RESET ( view_option_name [, ... ] )
```

Description

ALTER VIEW changes various auxiliary properties of a view. (If you want to modify the view's defining query, use CREATE OR REPLACE VIEW.)

You must own the view to use ALTER VIEW. To change a view's schema, you must also have CREATE privilege on the new schema. To alter the owner, you must be able to SET ROLE to the new owning role, and that role must have CREATE privilege on the view's schema. (These restrictions enforce that altering the owner doesn't do anything you couldn't do by dropping and recreating the view. However, a superuser can alter ownership of any view anyway.)

Parameters

name

The name (optionally schema-qualified) of an existing view.

column_name

Name of an existing column.

new_column_name

New name for an existing column.

IF EXISTS

Do not throw an error if the view does not exist. A notice is issued in this case.

SET/DROP DEFAULT

These forms set or remove the default value for a column. A view column's default value is substituted into any INSERT or UPDATE command whose target is the view, before applying any rules or triggers for the view. The view's default will therefore take precedence over any default values from underlying relations.

new_owner

The user name of the new owner of the view.

new_name

The new name for the view.

new_schema

The new schema for the view.

```
SET ( view_option_name [= view_option_value] [, ... ] )
RESET ( view_option_name [, ... ] )
```

Sets or resets a view option. Currently supported options are:

`check_option` (enum)

Changes the check option of the view. The value must be `local` or `cascaded`.

`security_barrier` (boolean)

Changes the security-barrier property of the view. The value must be a Boolean value, such as `true` or `false`.

`security_invoker` (boolean)

Changes the security-invoker property of the view. The value must be a Boolean value, such as `true` or `false`.

Notes

For historical reasons, `ALTER TABLE` can be used with views too; but the only variants of `ALTER TABLE` that are allowed with views are equivalent to the ones shown above.

Examples

To rename the view `foo` to `bar`:

```
ALTER VIEW foo RENAME TO bar;
```

To attach a default column value to an updatable view:

```
CREATE TABLE base_table (id int, ts timestamptz);
CREATE VIEW a_view AS SELECT * FROM base_table;
ALTER VIEW a_view ALTER COLUMN ts SET DEFAULT now();
INSERT INTO base_table(id) VALUES(1); -- ts will receive a NULL
INSERT INTO a_view(id) VALUES(2); -- ts will receive the current time
```

Compatibility

`ALTER VIEW` is a Postgres Pro extension of the SQL standard.

See Also

[CREATE VIEW](#), [DROP VIEW](#)

ANALYZE

ANALYZE — collect statistics about a database

Synopsis

```
ANALYZE [ ( option [, ...] ) ] [ table_and_columns [, ...] ]
```

where *option* can be one of:

```
VERBOSE [ boolean ]
SKIP_LOCKED [ boolean ]
BUFFER_USAGE_LIMIT size
```

and *table_and_columns* is:

```
table_name [ ( column_name [, ...] ) ]
```

Description

ANALYZE collects statistics about the contents of tables in the database, and stores the results in the `pg_statistic` system catalog. Subsequently, the query planner uses these statistics to help determine the most efficient execution plans for queries.

Without a *table_and_columns* list, ANALYZE processes every table and materialized view in the current database that the current user has permission to analyze. With a list, ANALYZE processes only those table(s). It is further possible to give a list of column names for a table, in which case only the statistics for those columns are collected.

For information about using ANALYZE in distributed systems, see [Analyzing and Vacuuming in Distributed System](#).

Parameters

VERBOSE

Enables display of progress messages.

SKIP_LOCKED

Specifies that ANALYZE should not wait for any conflicting locks to be released when beginning work on a relation: if a relation cannot be locked immediately without waiting, the relation is skipped. Note that even with this option, ANALYZE may still block when opening the relation's indexes or when acquiring sample rows from partitions, table inheritance children, and some types of foreign tables. Also, while ANALYZE ordinarily processes all partitions of specified partitioned tables, this option will cause ANALYZE to skip all partitions if there is a conflicting lock on the partitioned table.

BUFFER_USAGE_LIMIT

Specifies the [Buffer Access Strategy](#) ring buffer size for ANALYZE. This size is used to calculate the number of shared buffers which will be reused as part of this strategy. 0 disables use of a Buffer Access Strategy. When this option is not specified, ANALYZE uses the value from [vacuum_buffer_usage_limit](#). Higher settings can allow ANALYZE to run more quickly, but having too large a setting may cause too many other useful pages to be evicted from shared buffers. The minimum value is 128 kB and the maximum value is 16 GB.

boolean

Specifies whether the selected option should be turned on or off. You can write TRUE, ON, or 1 to enable the option, and FALSE, OFF, or 0 to disable it. The *boolean* value can also be omitted, in which case TRUE is assumed.

size

Specifies an amount of memory in kilobytes. Sizes may also be specified as a string containing the numerical size followed by any one of the following memory units: B (bytes), kB (kilobytes), MB (megabytes), GB (gigabytes), or TB (terabytes).

table_name

The name (possibly schema-qualified) of a specific table to analyze. If omitted, all regular tables, partitioned tables, and materialized views in the current database are analyzed (but not foreign tables). If the specified table is a partitioned table, both the inheritance statistics of the partitioned table as a whole and statistics of the individual partitions are updated.

column_name

The name of a specific column to analyze. Defaults to all columns.

Outputs

When `VERBOSE` is specified, `ANALYZE` emits progress messages to indicate which table is currently being processed. Various statistics about the tables are printed as well.

Notes

To analyze a table, one must ordinarily have the `MAINTAIN` privilege on the table. However, database owners are allowed to analyze all tables in their databases, except shared catalogs. `ANALYZE` will skip over any tables that the calling user does not have permission to analyze.

Foreign tables are analyzed only when explicitly selected. Not all foreign data wrappers support `ANALYZE`. If the table's wrapper does not support `ANALYZE`, the command prints a warning and does nothing.

In the default Postgres Pro configuration, the autovacuum daemon (see [Section 24.1.6](#)) takes care of automatic analyzing of tables when they are first loaded with data, and as they change throughout regular operation. When autovacuum is disabled, it is a good idea to run `ANALYZE` periodically, or just after making major changes in the contents of a table. Accurate statistics will help the planner to choose the most appropriate query plan, and thereby improve the speed of query processing. A common strategy for read-mostly databases is to run `VACUUM` and `ANALYZE` once a day during a low-usage time of day. (This will not be sufficient if there is heavy update activity.)

While `ANALYZE` is running, the [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`.

`ANALYZE` requires only a read lock on the target table, so it can run in parallel with other non-DDL activity on the table.

The statistics collected by `ANALYZE` usually include a list of some of the most common values in each column and a histogram showing the approximate data distribution in each column. One or both of these can be omitted if `ANALYZE` deems them uninteresting (for example, in a unique-key column, there are no common values) or if the column data type does not support the appropriate operators. There is more information about the statistics in [Chapter 24](#).

For large tables, `ANALYZE` takes a random sample of the table contents, rather than examining every row. This allows even very large tables to be analyzed in a small amount of time. Note, however, that the statistics are only approximate, and will change slightly each time `ANALYZE` is run, even if the actual table contents did not change. This might result in small changes in the planner's estimated costs shown by `EXPLAIN`. In rare situations, this non-determinism will cause the planner's choices of query plans to change after `ANALYZE` is run. To avoid this, raise the amount of statistics collected by `ANALYZE`, as described below.

The extent of analysis can be controlled by adjusting the [default_statistics_target](#) configuration variable, or on a column-by-column basis by setting the per-column statistics target with `ALTER TABLE ... ALTER COLUMN ... SET STATISTICS`. The target value sets the maximum number of entries in the most-

common-value list and the maximum number of bins in the histogram. The default target value is 100, but this can be adjusted up or down to trade off accuracy of planner estimates against the time taken for `ANALYZE` and the amount of space occupied in `pg_statistic`. In particular, setting the statistics target to zero disables collection of statistics for that column. It might be useful to do that for columns that are never used as part of the `WHERE`, `GROUP BY`, or `ORDER BY` clauses of queries, since the planner will have no use for statistics on such columns.

The largest statistics target among the columns being analyzed determines the number of table rows sampled to prepare the statistics. Increasing the target causes a proportional increase in the time and space needed to do `ANALYZE`.

One of the values estimated by `ANALYZE` is the number of distinct values that appear in each column. Because only a subset of the rows are examined, this estimate can sometimes be quite inaccurate, even with the largest possible statistics target. If this inaccuracy leads to bad query plans, a more accurate value can be determined manually and then installed with `ALTER TABLE ... ALTER COLUMN ... SET (n_distinct = ...)`.

If the table being analyzed has inheritance children, `ANALYZE` gathers two sets of statistics: one on the rows of the parent table only, and a second including rows of both the parent table and all of its children. This second set of statistics is needed when planning queries that process the inheritance tree as a whole. The child tables themselves are not individually analyzed in this case. The autovacuum daemon, however, will only consider inserts or updates on the parent table itself when deciding whether to trigger an automatic analyze for that table. If that table is rarely inserted into or updated, the inheritance statistics will not be up to date unless you run `ANALYZE` manually.

For partitioned tables, `ANALYZE` gathers statistics by sampling rows from all partitions; in addition, it will recurse into each partition and update its statistics. Each leaf partition is analyzed only once, even with multi-level partitioning. No statistics are collected for only the parent table (without data from its partitions), because with partitioning it's guaranteed to be empty.

The autovacuum daemon does not process partitioned tables, nor does it process inheritance parents if only the children are ever modified. It is usually necessary to periodically run a manual `ANALYZE` to keep the statistics of the table hierarchy up to date.

If any child tables or partitions are foreign tables whose foreign data wrappers do not support `ANALYZE`, those tables are ignored while gathering inheritance statistics.

If the table being analyzed is completely empty, `ANALYZE` will not record new statistics for that table. Any existing statistics will be retained.

Each backend running `ANALYZE` will report its progress in the `pg_stat_progress_analyze` view. See [Section 27.4.1](#) for details.

Compatibility

There is no `ANALYZE` statement in the SQL standard.

The following syntax was used before Postgres Pro version 11 and is still supported:

```
ANALYZE [ VERBOSE ] [ table_and_columns [, ...] ]
```

See Also

[VACUUM](#), [vacuumdb](#), [Section 19.4.4](#), [Section 24.1.6](#), [Section 27.4.1](#)

BEGIN

BEGIN — start a transaction block

Synopsis

```
BEGIN [ AUTONOMOUS ] [ WORK | TRANSACTION ] [ transaction_mode [, ...] ]
```

where *transaction_mode* is one of:

```
ISOLATION LEVEL { SERIALIZABLE | REPEATABLE READ | READ COMMITTED | READ  
UNCOMMITTED }  
READ WRITE | READ ONLY  
[ NOT ] DEFERRABLE
```

Description

BEGIN initiates a transaction block, that is, all statements after a BEGIN command will be executed in a single transaction until an explicit COMMIT or ROLLBACK is given. By default (without BEGIN), Postgres Pro executes transactions in “autocommit” mode, that is, each statement is executed in its own transaction and a commit is implicitly performed at the end of the statement (if execution was successful, otherwise a rollback is done).

Statements are executed more quickly in a transaction block, because transaction start/commit requires significant CPU and disk activity. Execution of multiple statements inside a transaction is also useful to ensure consistency when making several related changes: other sessions will be unable to see the intermediate states wherein not all the related updates have been done.

If the isolation level, read/write mode, or deferrable mode is specified, the new transaction has those characteristics, as if SET TRANSACTION was executed.

If you specify an optional AUTONOMOUS keyword, an autonomous transaction is started. Autonomous transactions can be started only within another transaction. For details, see [Chapter 16](#).

Parameters

AUTONOMOUS

Starts an autonomous transaction within another transaction.

WORK

TRANSACTION

Optional key words. They have no effect.

Refer to [SET TRANSACTION](#) for information on the meaning of the other parameters to this statement.

Notes

[START TRANSACTION](#) has the same functionality as BEGIN.

Use COMMIT or ROLLBACK to terminate a transaction block.

Issuing BEGIN when already inside a transaction block will provoke a warning message. The state of the transaction is not affected. To nest transactions within a transaction block, use savepoints (see [SAVEPOINT](#)).

For reasons of backwards compatibility, the commas between successive *transaction_modes* can be omitted.

Examples

To begin a transaction block:

```
BEGIN;
```

Compatibility

BEGIN is a Postgres Pro language extension. It is equivalent to the SQL-standard command [START TRANSACTION](#), whose reference page contains additional compatibility information.

The DEFERRABLE *transaction_mode* is a Postgres Pro language extension.

Incidentally, the BEGIN key word is used for a different purpose in embedded SQL. You are advised to be careful about the transaction semantics when porting database applications.

See Also

[COMMIT](#), [ROLLBACK](#), [START TRANSACTION](#), [SAVEPOINT](#)

CALL

CALL — invoke a procedure

Synopsis

```
CALL name ( [ argument ] [, ...] )
```

Description

CALL executes a procedure.

If the procedure has any output parameters, then a result row will be returned, containing the values of those parameters.

Parameters

name

The name (optionally schema-qualified) of the procedure.

argument

An argument expression for the procedure call.

Arguments can include parameter names, using the syntax *name* => *value*. This works the same as in ordinary function calls; see [Section 4.3](#) for details.

Arguments must be supplied for all procedure parameters that lack defaults, including OUT parameters. However, arguments matching OUT parameters are not evaluated, so it's customary to just write NULL for them. (Writing something else for an OUT parameter might cause compatibility problems with future Postgres Pro versions.)

Notes

The user must have EXECUTE privilege on the procedure in order to be allowed to invoke it.

To call a function (not a procedure), use SELECT instead.

If CALL is executed in a transaction block, then the called procedure cannot execute transaction control statements. Transaction control statements are only allowed if CALL is executed in its own transaction.

PL/pgSQL handles output parameters in CALL commands differently; see [Section 44.6.3](#).

Examples

```
CALL do_db_maintenance();
```

Compatibility

CALL conforms to the SQL standard, except for the handling of output parameters. The standard says that users should write variables to receive the values of output parameters.

See Also

[CREATE PROCEDURE](#)

CHECKPOINT

CHECKPOINT — force a write-ahead log checkpoint

Synopsis

CHECKPOINT

Description

A checkpoint is a point in the write-ahead log sequence at which all data files have been updated to reflect the information in the log. All data files will be flushed to disk. Refer to [Section 28.5](#) for more details about what happens during a checkpoint.

The `CHECKPOINT` command forces an immediate checkpoint when the command is issued, without waiting for a regular checkpoint scheduled by the system (controlled by the settings in [Section 19.5.2](#)). `CHECKPOINT` is not intended for use during normal operation.

If executed during recovery, the `CHECKPOINT` command will force a restartpoint (see [Section 28.5](#)) rather than writing a new checkpoint.

Only superusers or users with the privileges of the `pg_checkpoint` role can call `CHECKPOINT`.

Compatibility

The `CHECKPOINT` command is a Postgres Pro language extension.

CLOSE

CLOSE — close a cursor

Synopsis

```
CLOSE { name | ALL }
```

Description

CLOSE frees the resources associated with an open cursor. After the cursor is closed, no subsequent operations are allowed on it. A cursor should be closed when it is no longer needed.

Every non-holdable open cursor is implicitly closed when a transaction is terminated by COMMIT or ROLLBACK. A holdable cursor is implicitly closed if the transaction that created it aborts via ROLLBACK. If the creating transaction successfully commits, the holdable cursor remains open until an explicit CLOSE is executed, or the client disconnects.

Parameters

name

The name of an open cursor to close.

ALL

Close all open cursors.

Notes

Postgres Pro does not have an explicit OPEN cursor statement; a cursor is considered open when it is declared. Use the DECLARE statement to declare a cursor.

You can see all available cursors by querying the [pg_cursors](#) system view.

If a cursor is closed after a savepoint which is later rolled back, the CLOSE is not rolled back; that is, the cursor remains closed.

Examples

Close the cursor `liahona`:

```
CLOSE liahona;
```

Compatibility

CLOSE is fully conforming with the SQL standard. CLOSE ALL is a Postgres Pro extension.

See Also

[DECLARE](#), [FETCH](#), [MOVE](#)

CLUSTER

CLUSTER — cluster a table according to an index

Synopsis

```
CLUSTER [ ( option [, ...] ) ] [ table_name [ USING index_name ] ]
```

where *option* can be one of:

```
VERBOSE [ boolean ]
```

Description

CLUSTER instructs Postgres Pro to cluster the table specified by *table_name* based on the index specified by *index_name*. The index must already have been defined on *table_name*.

When a table is clustered, it is physically reordered based on the index information. Clustering is a one-time operation: when the table is subsequently updated, the changes are not clustered. That is, no attempt is made to store new or updated rows according to their index order. (If one wishes, one can periodically recluster by issuing the command again. Also, setting the table's `fillfactor` storage parameter to less than 100% can aid in preserving cluster ordering during updates, since updated rows are kept on the same page if enough space is available there.)

When a table is clustered, Postgres Pro remembers which index it was clustered by. The form `CLUSTER table_name` reclusters the table using the same index as before. You can also use the `CLUSTER` or `SET WITHOUT CLUSTER` forms of `ALTER TABLE` to set the index to be used for future cluster operations, or to clear any previous setting.

CLUSTER without a *table_name* reclusters all the previously-clustered tables in the current database that the calling user has privileges for. This form of CLUSTER cannot be executed inside a transaction block.

When a table is being clustered, an `ACCESS EXCLUSIVE` lock is acquired on it. This prevents any other database operations (both reads and writes) from operating on the table until the `CLUSTER` is finished.

Parameters

table_name

The name (possibly schema-qualified) of a table.

index_name

The name of an index.

VERBOSE

Prints a progress report as each table is clustered.

boolean

Specifies whether the selected option should be turned on or off. You can write `TRUE`, `ON`, or `1` to enable the option, and `FALSE`, `OFF`, or `0` to disable it. The *boolean* value can also be omitted, in which case `TRUE` is assumed.

Notes

To cluster a table, one must have the `MAINTAIN` privilege on the table.

In cases where you are accessing single rows randomly within a table, the actual order of the data in the table is unimportant. However, if you tend to access some data more than others, and there is an index

that groups them together, you will benefit from using `CLUSTER`. If you are requesting a range of indexed values from a table, or a single indexed value that has multiple rows that match, `CLUSTER` will help because once the index identifies the table page for the first row that matches, all other rows that match are probably already on the same table page, and so you save disk accesses and speed up the query.

`CLUSTER` can re-sort the table using either an index scan on the specified index, or (if the index is a b-tree) a sequential scan followed by sorting. It will attempt to choose the method that will be faster, based on planner cost parameters and available statistical information.

While `CLUSTER` is running, the [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`.

When an index scan is used, a temporary copy of the table is created that contains the table data in the index order. Temporary copies of each index on the table are created as well. Therefore, you need free space on disk at least equal to the sum of the table size and the index sizes.

When a sequential scan and sort is used, a temporary sort file is also created, so that the peak temporary space requirement is as much as double the table size, plus the index sizes. This method is often faster than the index scan method, but if the disk space requirement is intolerable, you can disable this choice by temporarily setting [enable_sort](#) to `off`.

It is advisable to set [maintenance_work_mem](#) to a reasonably large value (but not more than the amount of RAM you can dedicate to the `CLUSTER` operation) before clustering.

Because the planner records statistics about the ordering of tables, it is advisable to run [ANALYZE](#) on the newly clustered table. Otherwise, the planner might make poor choices of query plans.

Because `CLUSTER` remembers which indexes are clustered, one can cluster the tables one wants clustered manually the first time, then set up a periodic maintenance script that executes `CLUSTER` without any parameters, so that the desired tables are periodically reclustered.

Each backend running `CLUSTER` will report its progress in the `pg_stat_progress_cluster` view. See [Section 27.4.2](#) for details.

Clustering a partitioned table clusters each of its partitions using the partition of the specified partitioned index. When clustering a partitioned table, the index may not be omitted. `CLUSTER` on a partitioned table cannot be executed inside a transaction block.

Examples

Cluster the table `employees` on the basis of its index `employees_ind`:

```
CLUSTER employees USING employees_ind;
```

Cluster the `employees` table using the same index that was used before:

```
CLUSTER employees;
```

Cluster all tables in the database that have previously been clustered:

```
CLUSTER;
```

Compatibility

There is no `CLUSTER` statement in the SQL standard.

The following syntax was used before PostgreSQL 17 and is still supported:

```
CLUSTER [ VERBOSE ] [ table_name [ USING index_name ] ]
```

The following syntax was used before PostgreSQL 8.3 and is still supported:

```
CLUSTER index_name ON table_name
```

See Also

[clusterdb](#), [Section 27.4.2](#)

COMMENT

COMMENT — define or change the comment of an object

Synopsis

```
COMMENT ON
{
    ACCESS METHOD object_name |
    AGGREGATE aggregate_name ( aggregate_signature ) |
    CAST (source_type AS target_type) |
    COLLATION object_name |
    COLUMN relation_name.column_name |
    CONSTRAINT constraint_name ON table_name |
    CONSTRAINT constraint_name ON DOMAIN domain_name |
    CONVERSION object_name |
    DATABASE object_name |
    DOMAIN object_name |
    EXTENSION object_name |
    EVENT TRIGGER object_name |
    FOREIGN DATA WRAPPER object_name |
    FOREIGN TABLE object_name |
    FUNCTION function_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    INDEX object_name |
    LARGE OBJECT large_object_oid |
    MATERIALIZED VIEW object_name |
    OPERATOR operator_name (left_type, right_type) |
    OPERATOR CLASS object_name USING index_method |
    OPERATOR FAMILY object_name USING index_method |
    POLICY policy_name ON table_name |
    [ PROCEDURAL ] LANGUAGE object_name |
    PROCEDURE procedure_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    PROFILE object_name |
    PUBLICATION object_name |
    ROLE object_name |
    ROUTINE routine_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    RULE rule_name ON table_name |
    SCHEMA object_name |
    SEQUENCE object_name |
    SERVER object_name |
    STATISTICS object_name |
    SUBSCRIPTION object_name |
    TABLE object_name |
    TABLESPACE object_name |
    TEXT SEARCH CONFIGURATION object_name |
    TEXT SEARCH DICTIONARY object_name |
    TEXT SEARCH PARSER object_name |
    TEXT SEARCH TEMPLATE object_name |
    TRANSFORM FOR type_name LANGUAGE lang_name |
    TRIGGER trigger_name ON table_name |
    TYPE object_name |
    VIEW object_name
} IS { string_literal | NULL }
```

where *aggregate_signature* is:

```
* |
```

```
[ argmode ] [ argname ] argtype [ , ... ] |  
[ [ argmode ] [ argname ] argtype [ , ... ] ] ORDER BY [ argmode ] [ argname ] argtype  
[ , ... ]
```

Description

COMMENT stores a comment about a database object.

Only one comment string is stored for each object, so to modify a comment, issue a new COMMENT command for the same object. To remove a comment, write NULL in place of the text string. Comments are automatically dropped when their object is dropped.

A SHARE UPDATE EXCLUSIVE lock is acquired on the object to be commented.

For most kinds of object, only the object's owner can set the comment. Roles don't have owners, so the rule for COMMENT ON ROLE is that you must be superuser to comment on a superuser role, or have the CREATEROLE privilege and have been granted ADMIN OPTION on the target role. Likewise, access methods don't have owners either; you must be superuser to comment on an access method. Of course, a superuser can comment on anything.

Comments can be viewed using psql's \d family of commands. Other user interfaces to retrieve comments can be built atop the same built-in functions that psql uses, namely obj_description, col_description, and shobj_description (see [Table 9.80](#)).

Parameters

object_name
relation_name.column_name
aggregate_name
constraint_name
function_name
operator_name
policy_name
procedure_name
routine_name
rule_name
trigger_name

The name of the object to be commented. Names of objects that reside in schemas (tables, functions, etc.) can be schema-qualified. When commenting on a column, *relation_name* must refer to a table, view, composite type, or foreign table.

table_name
domain_name

When creating a comment on a constraint, a trigger, a rule or a policy these parameters specify the name of the table or domain on which that object is defined.

source_type

The name of the source data type of the cast.

target_type

The name of the target data type of the cast.

argmode

The mode of a function, procedure, or aggregate argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN. Note that COMMENT does not actually pay any attention to OUT arguments, since only the input arguments are needed to determine the function's identity. So it is sufficient to list the IN, INOUT, and VARIADIC arguments.

argname

The name of a function, procedure, or aggregate argument. Note that `COMMENT` does not actually pay any attention to argument names, since only the argument data types are needed to determine the function's identity.

argtype

The data type of a function, procedure, or aggregate argument.

large_object_oid

The OID of the large object.

left_type

right_type

The data type(s) of the operator's arguments (optionally schema-qualified). Write `NONE` for the missing argument of a prefix operator.

`PROCEDURAL`

This is a noise word.

type_name

The name of the data type of the transform.

lang_name

The name of the language of the transform.

string_literal

The new comment contents, written as a string literal.

`NULL`

Write `NULL` to drop the comment.

Notes

There is presently no security mechanism for viewing comments: any user connected to a database can see all the comments for objects in that database. For shared objects such as databases, roles, and tablespaces, comments are stored globally so any user connected to any database in the cluster can see all the comments for shared objects. Therefore, don't put security-critical information in comments.

Examples

Attach a comment to the table `mytable`:

```
COMMENT ON TABLE mytable IS 'This is my table.';
```

Remove it again:

```
COMMENT ON TABLE mytable IS NULL;
```

Some more examples:

```
COMMENT ON ACCESS METHOD gin IS 'GIN index access method';
COMMENT ON AGGREGATE my_aggregate (double precision) IS 'Computes sample variance';
COMMENT ON CAST (text AS int4) IS 'Allow casts from text to int4';
COMMENT ON COLLATION "fr_CA" IS 'Canadian French';
COMMENT ON COLUMN my_table.my_column IS 'Employee ID number';
COMMENT ON CONVERSION my_conv IS 'Conversion to UTF8';
COMMENT ON CONSTRAINT bar_col_cons ON bar IS 'Constrains column col';
```

```
COMMENT ON CONSTRAINT dom_col_constr ON DOMAIN dom IS 'Constrains col of domain';
COMMENT ON DATABASE my_database IS 'Development Database';
COMMENT ON DOMAIN my_domain IS 'Email Address Domain';
COMMENT ON EVENT TRIGGER abort_ddl IS 'Aborts all DDL commands';
COMMENT ON EXTENSION hstore IS 'implements the hstore data type';
COMMENT ON FOREIGN DATA WRAPPER mywrapper IS 'my foreign data wrapper';
COMMENT ON FOREIGN TABLE my_foreign_table IS 'Employee Information in other database';
COMMENT ON FUNCTION my_function (timestamp) IS 'Returns Roman Numeral';
COMMENT ON INDEX my_index IS 'Enforces uniqueness on employee ID';
COMMENT ON LANGUAGE plpython IS 'Python support for stored procedures';
COMMENT ON LARGE OBJECT 346344 IS 'Planning document';
COMMENT ON MATERIALIZED VIEW my_matview IS 'Summary of order history';
COMMENT ON OPERATOR ^ (text, text) IS 'Performs intersection of two texts';
COMMENT ON OPERATOR - (NONE, integer) IS 'Unary minus';
COMMENT ON OPERATOR CLASS int4ops USING btree IS '4 byte integer operators for btrees';
COMMENT ON OPERATOR FAMILY integer_ops USING btree IS 'all integer operators for
  btrees';
COMMENT ON POLICY my_policy ON mytable IS 'Filter rows by users';
COMMENT ON PROCEDURE my_proc (integer, integer) IS 'Runs a report';
COMMENT ON PUBLICATION alltables IS 'Publishes all operations on all tables';
COMMENT ON ROLE my_role IS 'Administration group for finance tables';
COMMENT ON ROUTINE my_routine (integer, integer) IS 'Runs a routine (which is a
  function or procedure)';
COMMENT ON RULE my_rule ON my_table IS 'Logs updates of employee records';
COMMENT ON SCHEMA my_schema IS 'Departmental data';
COMMENT ON SEQUENCE my_sequence IS 'Used to generate primary keys';
COMMENT ON SERVER myserver IS 'my foreign server';
COMMENT ON STATISTICS my_statistics IS 'Improves planner row estimations';
COMMENT ON SUBSCRIPTION alltables IS 'Subscription for all operations on all tables';
COMMENT ON TABLE my_schema.my_table IS 'Employee Information';
COMMENT ON TABLESPACE my_tablespace IS 'Tablespace for indexes';
COMMENT ON TEXT SEARCH CONFIGURATION my_config IS 'Special word filtering';
COMMENT ON TEXT SEARCH DICTIONARY swedish IS 'Snowball stemmer for Swedish language';
COMMENT ON TEXT SEARCH PARSER my_parser IS 'Splits text into words';
COMMENT ON TEXT SEARCH TEMPLATE snowball IS 'Snowball stemmer';
COMMENT ON TRANSFORM FOR hstore LANGUAGE plpython3u IS 'Transform between hstore and
  Python dict';
COMMENT ON TRIGGER my_trigger ON my_table IS 'Used for RI';
COMMENT ON TYPE complex IS 'Complex number data type';
COMMENT ON VIEW my_view IS 'View of departmental costs';
```

Compatibility

There is no `COMMENT` command in the SQL standard.

COMMIT

COMMIT — commit the current transaction

Synopsis

```
COMMIT [ AUTONOMOUS ] [ WORK | TRANSACTION ] [ AND [ NO ] CHAIN ]
```

Description

COMMIT commits the current transaction. All changes made by the transaction become visible to others and are guaranteed to be durable if a crash occurs.

Parameters

AUTONOMOUS

Optional key word that can be used when committing an autonomous transaction. For details on autonomous transactions, see [Chapter 16](#).

WORK

TRANSACTION

Optional key words. They have no effect.

AND CHAIN

If AND CHAIN is specified, a new transaction is immediately started with the same transaction characteristics (see [SET TRANSACTION](#)) as the just finished one. Otherwise, no new transaction is started.

Notes

Use [ROLLBACK](#) to abort a transaction.

Issuing COMMIT when not inside a transaction does no harm, but it will provoke a warning message. COMMIT AND CHAIN when not inside a transaction is an error.

Examples

To commit the current transaction and make all changes permanent:

```
COMMIT;
```

Compatibility

The command COMMIT conforms to the SQL standard. The form COMMIT TRANSACTION is a Postgres Pro extension.

See Also

[BEGIN](#), [ROLLBACK](#)

COMMIT PREPARED

COMMIT PREPARED — commit a transaction that was earlier prepared for two-phase commit

Synopsis

```
COMMIT PREPARED transaction_id
```

Description

COMMIT PREPARED commits a transaction that is in prepared state.

Parameters

transaction_id

The transaction identifier of the transaction that is to be committed.

Notes

To commit a prepared transaction, you must be either the same user that executed the transaction originally, or a superuser. But you do not have to be in the same session that executed the transaction.

This command cannot be executed inside a transaction block. The prepared transaction is committed immediately.

All currently available prepared transactions are listed in the [pg_prepared_xacts](#) system view.

Examples

Commit the transaction identified by the transaction identifier `foobar`:

```
COMMIT PREPARED 'foobar';
```

Compatibility

COMMIT PREPARED is a Postgres Pro extension. It is intended for use by external transaction management systems, some of which are covered by standards (such as X/Open XA), but the SQL side of those systems is not standardized.

See Also

[PREPARE TRANSACTION](#), [ROLLBACK PREPARED](#)

COPY

COPY — copy data between a file and a table

Synopsis

```
COPY table_name [ ( column_name [, ...] ) ]  
  FROM { 'filename' | PROGRAM 'command' | STDIN }  
  [ [ WITH ] ( option [, ...] ) ]  
  [ WHERE condition ]  
  
COPY { table_name [ ( column_name [, ...] ) ] | ( query ) }  
  TO { 'filename' | PROGRAM 'command' | STDOUT }  
  [ [ WITH ] ( option [, ...] ) ]
```

where *option* can be one of:

```
FORMAT format_name  
FREEZE [ boolean ]  
DELIMITER 'delimiter_character'  
NULL 'null_string'  
DEFAULT 'default_string'  
HEADER [ boolean | MATCH ]  
QUOTE 'quote_character'  
ESCAPE 'escape_character'  
FORCE_QUOTE { ( column_name [, ...] ) | * }  
FORCE_NOT_NULL { ( column_name [, ...] ) | * }  
FORCE_NULL { ( column_name [, ...] ) | * }  
ON_ERROR error_action  
ENCODING 'encoding_name'  
LOG_VERBOSITY verbosity
```

Description

COPY moves data between Postgres Pro tables and standard file-system files. COPY TO copies the contents of a table *to* a file, while COPY FROM copies data *from* a file to a table (appending the data to whatever is in the table already). COPY TO can also copy the results of a SELECT query.

If a column list is specified, COPY TO copies only the data in the specified columns to the file. For COPY FROM, each field in the file is inserted, in order, into the specified column. Table columns not specified in the COPY FROM column list will receive their default values.

COPY with a file name instructs the Postgres Pro server to directly read from or write to a file. The file must be accessible by the Postgres Pro user (the user ID the server runs as) and the name must be specified from the viewpoint of the server. When PROGRAM is specified, the server executes the given command and reads from the standard output of the program, or writes to the standard input of the program. The command must be specified from the viewpoint of the server, and be executable by the Postgres Pro user. When STDIN or STDOUT is specified, data is transmitted via the connection between the client and the server.

Each backend running COPY will report its progress in the `pg_stat_progress_copy` view. See [Section 27.4.3](#) for details.

By default, COPY will fail if it encounters an error during processing. For use cases where a best-effort attempt at loading the entire file is desired, the `ON_ERROR` clause can be used to specify some other behavior.

Parameters

table_name

The name (optionally schema-qualified) of an existing table.

column_name

An optional list of columns to be copied. If no column list is specified, all columns of the table except generated columns will be copied.

query

A `SELECT`, `VALUES`, `INSERT`, `UPDATE`, `DELETE`, or `MERGE` command whose results are to be copied. Note that parentheses are required around the query.

For `INSERT`, `UPDATE`, `DELETE`, and `MERGE` queries a `RETURNING` clause must be provided, and the target relation must not have a conditional rule, nor an `ALSO` rule, nor an `INSTEAD` rule that expands to multiple statements.

filename

The path name of the input or output file. An input file name can be an absolute or relative path, but an output file name must be an absolute path. Windows users might need to use an `E''` string and double any backslashes used in the path name.

PROGRAM

A command to execute. In `COPY FROM`, the input is read from standard output of the command, and in `COPY TO`, the output is written to the standard input of the command.

Note that the command is invoked by the shell, so if you need to pass any arguments that come from an untrusted source, you must be careful to strip or escape any special characters that might have a special meaning for the shell. For security reasons, it is best to use a fixed command string, or at least avoid including any user input in it.

STDIN

Specifies that input comes from the client application.

STDOUT

Specifies that output goes to the client application.

boolean

Specifies whether the selected option should be turned on or off. You can write `TRUE`, `ON`, or `1` to enable the option, and `FALSE`, `OFF`, or `0` to disable it. The *boolean* value can also be omitted, in which case `TRUE` is assumed.

FORMAT

Selects the data format to be read or written: `text`, `csv` (Comma Separated Values), or `binary`. The default is `text`.

FREEZE

Requests copying the data with rows already frozen, just as they would be after running the `VACUUM FREEZE` command. This is intended as a performance option for initial data loading. Rows will be frozen only if the table being loaded has been created or truncated in the current subtransaction, there are no cursors open and there are no older snapshots held by this transaction. It is currently not possible to perform a `COPY FREEZE` on a partitioned table. This option is only allowed in `COPY FROM`.

Note that all other sessions will immediately be able to see the data once it has been successfully loaded. This violates the normal rules of MVCC visibility and users should be aware of the potential problems this might cause.

DELIMITER

Specifies the character that separates columns within each row (line) of the file. The default is a tab character in text format, a comma in CSV format. This must be a single one-byte character. This option is not allowed when using `binary` format.

NULL

Specifies the string that represents a null value. The default is `\N` (backslash-N) in text format, and an unquoted empty string in CSV format. You might prefer an empty string even in text format for cases where you don't want to distinguish nulls from empty strings. This option is not allowed when using `binary` format.

Note

When using `COPY FROM`, any data item that matches this string will be stored as a null value, so you should make sure that you use the same string as you used with `COPY TO`.

DEFAULT

Specifies the string that represents a default value. Each time the string is found in the input file, the default value of the corresponding column will be used. This option is allowed only in `COPY FROM`, and only when not using `binary` format.

HEADER

Specifies that the file contains a header line with the names of each column in the file. On output, the first line contains the column names from the table. On input, the first line is discarded when this option is set to `true` (or equivalent Boolean value). If this option is set to `MATCH`, the number and names of the columns in the header line must match the actual column names of the table, in order; otherwise an error is raised. This option is not allowed when using `binary` format. The `MATCH` option is only valid for `COPY FROM` commands.

QUOTE

Specifies the quoting character to be used when a data value is quoted. The default is double-quote. This must be a single one-byte character. This option is allowed only when using CSV format.

ESCAPE

Specifies the character that should appear before a data character that matches the `QUOTE` value. The default is the same as the `QUOTE` value (so that the quoting character is doubled if it appears in the data). This must be a single one-byte character. This option is allowed only when using CSV format.

FORCE_QUOTE

Forces quoting to be used for all non-NULL values in each specified column. NULL output is never quoted. If `*` is specified, non-NULL values will be quoted in all columns. This option is allowed only in `COPY TO`, and only when using CSV format.

FORCE_NOT_NULL

Do not match the specified columns' values against the null string. In the default case where the null string is empty, this means that empty values will be read as zero-length strings rather than nulls, even when they are not quoted. If `*` is specified, the option will be applied to all columns. This option is allowed only in `COPY FROM`, and only when using CSV format.

FORCE_NULL

Match the specified columns' values against the null string, even if it has been quoted, and if a match is found set the value to NULL. In the default case where the null string is empty, this converts a

quoted empty string into NULL. If `*` is specified, the option will be applied to all columns. This option is allowed only in `COPY FROM`, and only when using `CSV` format.

ON_ERROR

Specifies how to behave when encountering an error converting a column's input value into its data type. An *error_action* value of `stop` means fail the command, while `ignore` means discard the input row and continue with the next one. The default is `stop`.

The `ignore` option is applicable only for `COPY FROM` when the `FORMAT` is `text` or `csv`.

A `NOTICE` message containing the ignored row count is emitted at the end of the `COPY FROM` if at least one row was discarded. When `LOG_VERBOSITY` option is set to `verbose`, a `NOTICE` message containing the line of the input file and the column name whose input conversion has failed is emitted for each discarded row.

ENCODING

Specifies that the file is encoded in the *encoding_name*. If this option is omitted, the current client encoding is used. See the Notes below for more details.

LOG_VERBOSITY

Specify the amount of messages emitted by a `COPY` command: `default` or `verbose`. If `verbose` is specified, additional messages are emitted during processing.

This is currently used in `COPY FROM` command when `ON_ERROR` option is set to `ignore`.

WHERE

The optional `WHERE` clause has the general form

`WHERE condition`

where *condition* is any expression that evaluates to a result of type `boolean`. Any row that does not satisfy this condition will not be inserted to the table. A row satisfies the condition if it returns true when the actual row values are substituted for any variable references.

Currently, subqueries are not allowed in `WHERE` expressions, and the evaluation does not see any changes made by the `COPY` itself (this matters when the expression contains calls to `VOLATILE` functions).

Outputs

On successful completion, a `COPY` command returns a command tag of the form

`COPY count`

The *count* is the number of rows copied.

Note

psql will print this command tag only if the command was not `COPY ... TO STDOUT`, or the equivalent psql meta-command `\copy ... to stdout`. This is to prevent confusing the command tag with the data that was just printed.

Notes

`COPY TO` can be used only with plain tables, not views, and does not copy rows from child tables or child partitions. For example, `COPY table TO` copies the same rows as `SELECT * FROM ONLY table`. The syntax `COPY (SELECT * FROM table) TO ...` can be used to dump all of the rows in an inheritance hierarchy, partitioned table, or view.

`COPY FROM` can be used with plain, foreign, or partitioned tables or with views that have `INSTEAD OF INSERT` triggers.

You must have `select` privilege on the table whose values are read by `COPY TO`, and `insert` privilege on the table into which values are inserted by `COPY FROM`. It is sufficient to have column privileges on the column(s) listed in the command.

If row-level security is enabled for the table, the relevant `SELECT` policies will apply to `COPY table TO` statements. Currently, `COPY FROM` is not supported for tables with row-level security. Use equivalent `INSERT` statements instead.

Files named in a `COPY` command are read or written directly by the server, not by the client application. Therefore, they must reside on or be accessible to the database server machine, not the client. They must be accessible to and readable or writable by the Postgres Pro user (the user ID the server runs as), not the client. Similarly, the command specified with `PROGRAM` is executed directly by the server, not by the client application, must be executable by the Postgres Pro user. `COPY` naming a file or command is only allowed to database superusers or users who are granted one of the roles `pg_read_server_files`, `pg_write_server_files`, or `pg_execute_server_program`, since it allows reading or writing any file or running a program that the server has privileges to access.

Do not confuse `COPY` with the `psql` instruction `\copy`. `\copy` invokes `COPY FROM STDIN` or `COPY TO STDOUT`, and then fetches/stores the data in a file accessible to the `psql` client. Thus, file accessibility and access rights depend on the client rather than the server when `\copy` is used.

It is recommended that the file name used in `COPY` always be specified as an absolute path. This is enforced by the server in the case of `COPY TO`, but for `COPY FROM` you do have the option of reading from a file specified by a relative path. The path will be interpreted relative to the working directory of the server process (normally the cluster's data directory), not the client's working directory.

Executing a command with `PROGRAM` might be restricted by the operating system's access control mechanisms, such as SELinux.

`COPY FROM` will invoke any triggers and check constraints on the destination table. However, it will not invoke rules.

For identity columns, the `COPY FROM` command will always write the column values provided in the input data, like the `INSERT` option `OVERRIDING SYSTEM VALUE`.

`COPY` input and output is affected by `DateStyle`. To ensure portability to other Postgres Pro installations that might use non-default `DateStyle` settings, `DateStyle` should be set to `ISO` before using `COPY TO`. It is also a good idea to avoid dumping data with `IntervalStyle` set to `sql_standard`, because negative interval values might be misinterpreted by a server that has a different setting for `IntervalStyle`.

Input data is interpreted according to `ENCODING` option or the current client encoding, and output data is encoded in `ENCODING` or the current client encoding, even if the data does not pass through the client but is read from or written to a file directly by the server.

The `COPY FROM` command physically inserts input rows into the table as it progresses. If the command fails, these rows are left in a deleted state; these rows will not be visible, but still occupy disk space. This might amount to considerable wasted disk space if the failure happened well into a large copy operation. `VACUUM` should be used to recover the wasted space.

`FORCE_NULL` and `FORCE_NOT_NULL` can be used simultaneously on the same column. This results in converting quoted null strings to null values and unquoted null strings to empty strings.

Postgres Pro does not allow NUL bytes in data. If you are going to import such data using the `COPY FROM` command, you can specify an ASCII character in the [nul_byte_replacement_on_import](#) configuration parameter to replace NUL bytes on the fly.

File Formats

Text Format

When the `text` format is used, the data read or written is a text file with one line per table row. Columns in a row are separated by the delimiter character. The column values themselves are strings generated by the output function, or acceptable to the input function, of each attribute's data type. The specified null string is used in place of columns that are null. `COPY FROM` will raise an error if any line of the input file contains more or fewer columns than are expected.

End of data can be represented by a single line containing just backslash-period (`\.`). An end-of-data marker is not necessary when reading from a file, since the end of file serves perfectly well; it is needed only when copying data to or from client applications using pre-3.0 client protocol.

Backslash characters (`\`) can be used in the `COPY` data to quote data characters that might otherwise be taken as row or column delimiters. In particular, the following characters *must* be preceded by a backslash if they appear as part of a column value: backslash itself, newline, carriage return, and the current delimiter character.

The specified null string is sent by `COPY TO` without adding any backslashes; conversely, `COPY FROM` matches the input against the null string before removing backslashes. Therefore, a null string such as `\N` cannot be confused with the actual data value `\N` (which would be represented as `\\N`).

The following special backslash sequences are recognized by `COPY FROM`:

Sequence	Represents
<code>\b</code>	Backspace (ASCII 8)
<code>\f</code>	Form feed (ASCII 12)
<code>\n</code>	Newline (ASCII 10)
<code>\r</code>	Carriage return (ASCII 13)
<code>\t</code>	Tab (ASCII 9)
<code>\v</code>	Vertical tab (ASCII 11)
<code>\digits</code>	Backslash followed by one to three octal digits specifies the byte with that numeric code
<code>\xdigits</code>	Backslash <code>x</code> followed by one or two hex digits specifies the byte with that numeric code

Presently, `COPY TO` will never emit an octal or hex-digits backslash sequence, but it does use the other sequences listed above for those control characters.

Any other backslashed character that is not mentioned in the above table will be taken to represent itself. However, beware of adding backslashes unnecessarily, since that might accidentally produce a string matching the end-of-data marker (`\.`) or the null string (`\N` by default). These strings will be recognized before any other backslash processing is done.

It is strongly recommended that applications generating `COPY` data convert data newlines and carriage returns to the `\n` and `\r` sequences respectively. At present it is possible to represent a data carriage return by a backslash and carriage return, and to represent a data newline by a backslash and newline. However, these representations might not be accepted in future releases. They are also highly vulnerable to corruption if the `COPY` file is transferred across different machines (for example, from Unix to Windows or vice versa).

All backslash sequences are interpreted after encoding conversion. The bytes specified with the octal and hex-digit backslash sequences must form valid characters in the database encoding.

`COPY TO` will terminate each row with a Unix-style newline (`"\n"`). Servers running on Microsoft Windows instead output carriage return/newline (`"\r\n"`), but only for `COPY` to a server file; for consistency across

platforms, `COPY TO STDOUT` always sends “\n” regardless of server platform. `COPY FROM` can handle lines ending with newlines, carriage returns, or carriage return/newlines. To reduce the risk of error due to un-backslashed newlines or carriage returns that were meant as data, `COPY FROM` will complain if the line endings in the input are not all alike.

CSV Format

This format option is used for importing and exporting the Comma Separated Value (CSV) file format used by many other programs, such as spreadsheets. Instead of the escaping rules used by Postgres Pro's standard text format, it produces and recognizes the common CSV escaping mechanism.

The values in each record are separated by the `DELIMITER` character. If the value contains the delimiter character, the `QUOTE` character, the `NULL` string, a carriage return, or line feed character, then the whole value is prefixed and suffixed by the `QUOTE` character, and any occurrence within the value of a `QUOTE` character or the `ESCAPE` character is preceded by the escape character. You can also use `FORCE_QUOTE` to force quotes when outputting non-`NULL` values in specific columns.

The CSV format has no standard way to distinguish a `NULL` value from an empty string. Postgres Pro's `COPY` handles this by quoting. A `NULL` is output as the `NULL` parameter string and is not quoted, while a non-`NULL` value matching the `NULL` parameter string is quoted. For example, with the default settings, a `NULL` is written as an unquoted empty string, while an empty string data value is written with double quotes (“”). Reading values follows similar rules. You can use `FORCE_NOT_NULL` to prevent `NULL` input comparisons for specific columns. You can also use `FORCE_NULL` to convert quoted null string data values to `NULL`.

Because backslash is not a special character in the CSV format, `\.`, the end-of-data marker, could also appear as a data value. To avoid any misinterpretation, a `\.` data value appearing as a lone entry on a line is automatically quoted on output, and on input, if quoted, is not interpreted as the end-of-data marker. If you are loading a file created by another application that has a single unquoted column and might have a value of `\.`, you might need to quote that value in the input file.

Note

In CSV format, all characters are significant. A quoted value surrounded by white space, or any characters other than `DELIMITER`, will include those characters. This can cause errors if you import data from a system that pads CSV lines with white space out to some fixed width. If such a situation arises you might need to preprocess the CSV file to remove the trailing white space, before importing the data into Postgres Pro.

Note

CSV format will both recognize and produce CSV files with quoted values containing embedded carriage returns and line feeds. Thus the files are not strictly one line per table row like text-format files.

Note

Many programs produce strange and occasionally perverse CSV files, so the file format is more a convention than a standard. Thus you might encounter some files that cannot be imported using this mechanism, and `COPY` might produce files that other programs cannot process.

Binary Format

The `binary` format option causes all data to be stored/read as binary format rather than as text. It is somewhat faster than the text and CSV formats, but a binary-format file is less portable across machine

architectures and Postgres Pro versions. Also, the binary format is very data type specific; for example it will not work to output binary data from a `smallint` column and read it into an `integer` column, even though that would work fine in text format.

The `binary` file format consists of a file header, zero or more tuples containing the row data, and a file trailer. Headers and data are in network byte order.

Note

PostgreSQL releases before 7.4 used a different binary file format.

File Header

The file header consists of 15 bytes of fixed fields, followed by a variable-length header extension area. The fixed fields are:

Signature

11-byte sequence `PGCOPY\n\377\r\n\0` — note that the zero byte is a required part of the signature. (The signature is designed to allow easy identification of files that have been munged by a non-8-bit-clean transfer. This signature will be changed by end-of-line-translation filters, dropped zero bytes, dropped high bits, or parity changes.)

Flags field

32-bit integer bit mask to denote important aspects of the file format. Bits are numbered from 0 (LSB) to 31 (MSB). Note that this field is stored in network byte order (most significant byte first), as are all the integer fields used in the file format. Bits 16–31 are reserved to denote critical file format issues; a reader should abort if it finds an unexpected bit set in this range. Bits 0–15 are reserved to signal backwards-compatible format issues; a reader should simply ignore any unexpected bits set in this range. Currently only one flag bit is defined, and the rest must be zero:

Bit 16

If 1, OIDs are included in the data; if 0, not. Oid system columns are not supported in Postgres Pro anymore, but the format still contains the indicator.

Header extension area length

32-bit integer, length in bytes of remainder of header, not including self. Currently, this is zero, and the first tuple follows immediately. Future changes to the format might allow additional data to be present in the header. A reader should silently skip over any header extension data it does not know what to do with.

The header extension area is envisioned to contain a sequence of self-identifying chunks. The flags field is not intended to tell readers what is in the extension area. Specific design of header extension contents is left for a later release.

This design allows for both backwards-compatible header additions (add header extension chunks, or set low-order flag bits) and non-backwards-compatible changes (set high-order flag bits to signal such changes, and add supporting data to the extension area if needed).

Tuples

Each tuple begins with a 16-bit integer count of the number of fields in the tuple. (Presently, all tuples in a table will have the same count, but that might not always be true.) Then, repeated for each field in the tuple, there is a 32-bit length word followed by that many bytes of field data. (The length word does not include itself, and can be zero.) As a special case, -1 indicates a NULL field value. No value bytes follow in the NULL case.

There is no alignment padding or any other extra data between fields.

Presently, all data values in a binary-format file are assumed to be in binary format (format code one). It is anticipated that a future extension might add a header field that allows per-column format codes to be specified.

If OIDs are included in the file, the OID field immediately follows the field-count word. It is a normal field except that it's not included in the field-count. Note that oid system columns are not supported in current versions of Postgres Pro.

File Trailer

The file trailer consists of a 16-bit integer word containing -1. This is easily distinguished from a tuple's field-count word.

A reader should report an error if a field-count word is neither -1 nor the expected number of columns. This provides an extra check against somehow getting out of sync with the data.

Examples

The following example copies a table to the client using the vertical bar (|) as the field delimiter:

```
COPY country TO STDOUT (DELIMITER '|');
```

To copy data from a file into the `country` table:

```
COPY country FROM '/usr1/proj/bray/sql/country_data';
```

To copy into a file just the countries whose names start with 'A':

```
COPY (SELECT * FROM country WHERE country_name LIKE 'A%') TO '/usr1/proj/bray/sql/a_list_countries.copy';
```

To copy into a compressed file, you can pipe the output through an external compression program:

```
COPY country TO PROGRAM 'gzip > /usr1/proj/bray/sql/country_data.gz';
```

Here is a sample of data suitable for copying into a table from STDIN:

```
AF      AFGHANISTAN
AL      ALBANIA
DZ      ALGERIA
ZM      ZAMBIA
ZW      ZIMBABWE
```

Note that the white space on each line is actually a tab character.

The following is the same data, output in binary format. The data is shown after filtering through the Unix utility `od -c`. The table has three columns; the first has type `char(2)`, the second has type `text`, and the third has type `integer`. All the rows have a null value in the third column.

```
00000000  P   G   C   O   P   Y  \n 377  \r  \n  \0  \0  \0  \0  \0  \0
00000020  \0  \0  \0  \0  \0 003  \0  \0  \0 002  A   F  \0  \0  \0 013  A
00000040  F   G   H   A   N   I   S   T   A   N 377 377 377 377  \0 003
00000060  \0  \0  \0 002  A   L  \0  \0 007  A   L   B   A   N   I
00000100  A 377 377 377 377  \0 003  \0  \0  \0 002  D   Z  \0  \0  \0
00000120 007  A   L   G   E   R   I   A 377 377 377 377  \0 003  \0  \0
00000140  \0 002  Z   M  \0  \0  \0 006  Z   A   M   B   I   A 377 377
00000160 377 377  \0 003  \0  \0  \0 002  Z   W  \0  \0  \0  \b  Z   I
00000200  M   B   A   B   W   E 377 377 377 377 377 377
```

Compatibility

There is no `COPY` statement in the SQL standard.

The following syntax was used before PostgreSQL version 9.0 and is still supported:

```
COPY table_name [ ( column_name [, ...] ) ]
FROM { 'filename' | STDIN }
[ [ WITH ]
    [ BINARY ]
    [ DELIMITER [ AS ] 'delimiter_character' ]
    [ NULL [ AS ] 'null_string' ]
    [ CSV [ HEADER ]
        [ QUOTE [ AS ] 'quote_character' ]
        [ ESCAPE [ AS ] 'escape_character' ]
        [ FORCE NOT NULL column_name [, ...] ] ] ] ]

COPY { table_name [ ( column_name [, ...] ) ] | ( query ) }
TO { 'filename' | STDOUT }
[ [ WITH ]
    [ BINARY ]
    [ DELIMITER [ AS ] 'delimiter_character' ]
    [ NULL [ AS ] 'null_string' ]
    [ CSV [ HEADER ]
        [ QUOTE [ AS ] 'quote_character' ]
        [ ESCAPE [ AS ] 'escape_character' ]
        [ FORCE QUOTE { column_name [, ...] | * } ] ] ] ]
```

Note that in this syntax, **BINARY** and **CSV** are treated as independent keywords, not as arguments of a **FORMAT** option.

The following syntax was used before PostgreSQL version 7.3 and is still supported:

```
COPY [ BINARY ] table_name
FROM { 'filename' | STDIN }
[ [USING] DELIMITERS 'delimiter_character' ]
[ WITH NULL AS 'null_string' ]

COPY [ BINARY ] table_name
TO { 'filename' | STDOUT }
[ [USING] DELIMITERS 'delimiter_character' ]
[ WITH NULL AS 'null_string' ]
```

See Also

[Section 27.4.3](#)

CREATE ACCESS METHOD

CREATE ACCESS METHOD — define a new access method

Synopsis

```
CREATE ACCESS METHOD name
    TYPE access_method_type
    HANDLER handler_function
```

Description

CREATE ACCESS METHOD creates a new access method.

The access method name must be unique within the database.

Only superusers can define new access methods.

Parameters

name

The name of the access method to be created.

access_method_type

This clause specifies the type of access method to define. Only TABLE and INDEX are supported at present.

handler_function

handler_function is the name (possibly schema-qualified) of a previously registered function that represents the access method. The handler function must be declared to take a single argument of type `internal`, and its return type depends on the type of access method; for TABLE access methods, it must be `table_am_handler` and for INDEX access methods, it must be `index_am_handler`. The C-level API that the handler function must implement varies depending on the type of access method. The table access method API is described in [Chapter 62](#) and the index access method API is described in [Chapter 63](#).

Examples

Create an index access method `heptree` with handler function `heptree_handler`:

```
CREATE ACCESS METHOD heptree TYPE INDEX HANDLER heptree_handler;
```

Compatibility

CREATE ACCESS METHOD is a Postgres Pro extension.

See Also

[DROP ACCESS METHOD](#), [CREATE OPERATOR CLASS](#), [CREATE OPERATOR FAMILY](#)

CREATE AGGREGATE

CREATE AGGREGATE — define a new aggregate function

Synopsis

```
CREATE [ OR REPLACE ] AGGREGATE name ( [ argmode ] [ argname ] arg_data_type
[ , ... ] ) (
    SFUNC = sfunc,
    STYPE = state_data_type
    [ , SSPACE = state_data_size ]
    [ , FINALFUNC = ffunc ]
    [ , FINALFUNC_EXTRA ]
    [ , FINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE } ]
    [ , COMBINEFUNC = combinefunc ]
    [ , SERIALFUNC = serialfunc ]
    [ , DESERIALFUNC = deserialfunc ]
    [ , INITCOND = initial_condition ]
    [ , MSFUNC = msfunc ]
    [ , MINVFUNC = minvfunc ]
    [ , MSTYPE = mstate_data_type ]
    [ , MSSPACE = mstate_data_size ]
    [ , MFINALFUNC = mffunc ]
    [ , MFINALFUNC_EXTRA ]
    [ , MFINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE } ]
    [ , MINITCOND = minitial_condition ]
    [ , SORTOP = sort_operator ]
    [ , PARALLEL = { SAFE | RESTRICTED | UNSAFE } ]
)
```

```
CREATE [ OR REPLACE ] AGGREGATE name ( [ [ argmode ] [ argname ] arg_data_type
[ , ... ] ]
                                ORDER BY [ argmode ] [ argname ] arg_data_type [ , ... ] ) (
    SFUNC = sfunc,
    STYPE = state_data_type
    [ , SSPACE = state_data_size ]
    [ , FINALFUNC = ffunc ]
    [ , FINALFUNC_EXTRA ]
    [ , FINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE } ]
    [ , INITCOND = initial_condition ]
    [ , PARALLEL = { SAFE | RESTRICTED | UNSAFE } ]
    [ , HYPOTHETICAL ]
)
```

or the old syntax

```
CREATE [ OR REPLACE ] AGGREGATE name (
    BASETYPE = base_type,
    SFUNC = sfunc,
    STYPE = state_data_type
    [ , SSPACE = state_data_size ]
    [ , FINALFUNC = ffunc ]
    [ , FINALFUNC_EXTRA ]
    [ , FINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE } ]
    [ , COMBINEFUNC = combinefunc ]
    [ , SERIALFUNC = serialfunc ]
)
```



```
[ , DESERIALFUNC = deserialfunc ]
[ , INITCOND = initial_condition ]
[ , MSFUNC = msfunc ]
[ , MINVFUNC = minvfunc ]
[ , MSTYPE = mstate_data_type ]
[ , MSSPACE = mstate_data_size ]
[ , MFINALFUNC = mffunc ]
[ , MFINALFUNC_EXTRA ]
[ , MFINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE } ]
[ , MINITCOND = minitial_condition ]
[ , SORTOP = sort_operator ]
)
```

Description

`CREATE AGGREGATE` defines a new aggregate function. `CREATE OR REPLACE AGGREGATE` will either define a new aggregate function or replace an existing definition. Some basic and commonly-used aggregate functions are included with the distribution; they are documented in [Section 9.21](#). If one defines new types or needs an aggregate function not already provided, then `CREATE AGGREGATE` can be used to provide the desired features.

When replacing an existing definition, the argument types, result type, and number of direct arguments may not be changed. Also, the new definition must be of the same kind (ordinary aggregate, ordered-set aggregate, or hypothetical-set aggregate) as the old one.

If a schema name is given (for example, `CREATE AGGREGATE myschema.myagg ...`) then the aggregate function is created in the specified schema. Otherwise it is created in the current schema.

An aggregate function is identified by its name and input data type(s). Two aggregates in the same schema can have the same name if they operate on different input types. The name and input data type(s) of an aggregate must also be distinct from the name and input data type(s) of every ordinary function in the same schema. This behavior is identical to overloading of ordinary function names (see [CREATE FUNCTION](#)).

A simple aggregate function is made from one or two ordinary functions: a state transition function *sfunc*, and an optional final calculation function *ffunc*. These are used as follows:

```
sfunc( internal-state, next-data-values ) ---> next-internal-state
ffunc( internal-state ) ---> aggregate-value
```

Postgres Pro creates a temporary variable of data type *stype* to hold the current internal state of the aggregate. At each input row, the aggregate argument value(s) are calculated and the state transition function is invoked with the current state value and the new argument value(s) to calculate a new internal state value. After all the rows have been processed, the final function is invoked once to calculate the aggregate's return value. If there is no final function then the ending state value is returned as-is.

An aggregate function can provide an initial condition, that is, an initial value for the internal state value. This is specified and stored in the database as a value of type `text`, but it must be a valid external representation of a constant of the state value data type. If it is not supplied then the state value starts out null.

If the state transition function is declared “strict”, then it cannot be called with null inputs. With such a transition function, aggregate execution behaves as follows. Rows with any null input values are ignored (the function is not called and the previous state value is retained). If the initial state value is null, then at the first row with all-nonnull input values, the first argument value replaces the state value, and the transition function is invoked at each subsequent row with all-nonnull input values. This is handy for implementing aggregates like `max`. Note that this behavior is only available when *state_data_type* is the same as the first *arg_data_type*. When these types are different, you must supply a nonnull initial condition or use a nonstrict transition function.

If the state transition function is not strict, then it will be called unconditionally at each input row, and must deal with null inputs and null state values for itself. This allows the aggregate author to have full control over the aggregate's handling of null values.

If the final function is declared “strict”, then it will not be called when the ending state value is null; instead a null result will be returned automatically. (Of course this is just the normal behavior of strict functions.) In any case the final function has the option of returning a null value. For example, the final function for `avg` returns null when it sees there were zero input rows.

Sometimes it is useful to declare the final function as taking not just the state value, but extra parameters corresponding to the aggregate's input values. The main reason for doing this is if the final function is polymorphic and the state value's data type would be inadequate to pin down the result type. These extra parameters are always passed as NULL (and so the final function must not be strict when the `FINALFUNC_EXTRA` option is used), but nonetheless they are valid parameters. The final function could for example make use of `get_fn_expr_argtype` to identify the actual argument type in the current call.

An aggregate can optionally support *moving-aggregate mode*, as described in [Section 39.12.1](#). This requires specifying the `MSFUNC`, `MINVFUNC`, and `MSTYPE` parameters, and optionally the `MSSPACE`, `MFINALFUNC`, `MFINALFUNC_EXTRA`, `MFINALFUNC_MODIFY`, and `MINITCOND` parameters. Except for `MINVFUNC`, these parameters work like the corresponding simple-aggregate parameters without `M`; they define a separate implementation of the aggregate that includes an inverse transition function.

The syntax with `ORDER BY` in the parameter list creates a special type of aggregate called an *ordered-set aggregate*; or if `HYPOTHETICAL` is specified, then a *hypothetical-set aggregate* is created. These aggregates operate over groups of sorted values in order-dependent ways, so that specification of an input sort order is an essential part of a call. Also, they can have *direct* arguments, which are arguments that are evaluated only once per aggregation rather than once per input row. Hypothetical-set aggregates are a subclass of ordered-set aggregates in which some of the direct arguments are required to match, in number and data types, the aggregated argument columns. This allows the values of those direct arguments to be added to the collection of aggregate-input rows as an additional “hypothetical” row.

An aggregate can optionally support *partial aggregation*, as described in [Section 39.12.4](#). This requires specifying the `COMBINEFUNC` parameter. If the `state_data_type` is `internal`, it's usually also appropriate to provide the `SERIALFUNC` and `DESERIALFUNC` parameters so that parallel aggregation is possible. Note that the aggregate must also be marked `PARALLEL SAFE` to enable parallel aggregation.

Aggregates that behave like `MIN` or `MAX` can sometimes be optimized by looking into an index instead of scanning every input row. If this aggregate can be so optimized, indicate it by specifying a *sort operator*. The basic requirement is that the aggregate must yield the first element in the sort ordering induced by the operator; in other words:

```
SELECT agg(col) FROM tab;
```

must be equivalent to:

```
SELECT col FROM tab ORDER BY col USING sortop LIMIT 1;
```

Further assumptions are that the aggregate ignores null inputs, and that it delivers a null result if and only if there were no non-null inputs. Ordinarily, a data type's `<` operator is the proper sort operator for `MIN`, and `>` is the proper sort operator for `MAX`. Note that the optimization will never actually take effect unless the specified operator is the “less than” or “greater than” strategy member of a B-tree index operator class.

To be able to create an aggregate function, you must have `USAGE` privilege on the argument types, the state type(s), and the return type, as well as `EXECUTE` privilege on the supporting functions.

Parameters

name

The name (optionally schema-qualified) of the aggregate function to create.

argmode

The mode of an argument: `IN` or `VARIADIC`. (Aggregate functions do not support `OUT` arguments.) If omitted, the default is `IN`. Only the last argument can be marked `VARIADIC`.

argname

The name of an argument. This is currently only useful for documentation purposes. If omitted, the argument has no name.

arg_data_type

An input data type on which this aggregate function operates. To create a zero-argument aggregate function, write `*` in place of the list of argument specifications. (An example of such an aggregate is `count(*)`.)

base_type

In the old syntax for `CREATE AGGREGATE`, the input data type is specified by a `basetype` parameter rather than being written next to the aggregate name. Note that this syntax allows only one input parameter. To define a zero-argument aggregate function with this syntax, specify the `basetype` as `"ANY"` (not `*`). Ordered-set aggregates cannot be defined with the old syntax.

sfunc

The name of the state transition function to be called for each input row. For a normal N -argument aggregate function, the *sfunc* must take $N+1$ arguments, the first being of type *state_data_type* and the rest matching the declared input data type(s) of the aggregate. The function must return a value of type *state_data_type*. This function takes the current state value and the current input data value(s), and returns the next state value.

For ordered-set (including hypothetical-set) aggregates, the state transition function receives only the current state value and the aggregated arguments, not the direct arguments. Otherwise it is the same.

state_data_type

The data type for the aggregate's state value.

state_data_size

The approximate average size (in bytes) of the aggregate's state value. If this parameter is omitted or is zero, a default estimate is used based on the *state_data_type*. The planner uses this value to estimate the memory required for a grouped aggregate query.

ffunc

The name of the final function called to compute the aggregate's result after all input rows have been traversed. For a normal aggregate, this function must take a single argument of type *state_data_type*. The return data type of the aggregate is defined as the return type of this function. If *ffunc* is not specified, then the ending state value is used as the aggregate's result, and the return type is *state_data_type*.

For ordered-set (including hypothetical-set) aggregates, the final function receives not only the final state value, but also the values of all the direct arguments.

If `FINALFUNC_EXTRA` is specified, then in addition to the final state value and any direct arguments, the final function receives extra `NULL` values corresponding to the aggregate's regular (aggregated) arguments. This is mainly useful to allow correct resolution of the aggregate result type when a polymorphic aggregate is being defined.

`FINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE }`

This option specifies whether the final function is a pure function that does not modify its arguments. `READ_ONLY` indicates it does not; the other two values indicate that it may change the transition state

value. See [Notes](#) below for more detail. The default is `READ_ONLY`, except for ordered-set aggregates, for which the default is `READ_WRITE`.

combinefunc

The *combinefunc* function may optionally be specified to allow the aggregate function to support partial aggregation. If provided, the *combinefunc* must combine two *state_data_type* values, each containing the result of aggregation over some subset of the input values, to produce a new *state_data_type* that represents the result of aggregating over both sets of inputs. This function can be thought of as an *sfunc*, where instead of acting upon an individual input row and adding it to the running aggregate state, it adds another aggregate state to the running state.

The *combinefunc* must be declared as taking two arguments of the *state_data_type* and returning a value of the *state_data_type*. Optionally this function may be “strict”. In this case the function will not be called when either of the input states are null; the other state will be taken as the correct result.

For aggregate functions whose *state_data_type* is *internal*, the *combinefunc* must not be strict. In this case the *combinefunc* must ensure that null states are handled correctly and that the state being returned is properly stored in the aggregate memory context.

serialfunc

An aggregate function whose *state_data_type* is *internal* can participate in parallel aggregation only if it has a *serialfunc* function, which must serialize the aggregate state into a *bytea* value for transmission to another process. This function must take a single argument of type *internal* and return type *bytea*. A corresponding *deserialfunc* is also required.

deserialfunc

Deserialize a previously serialized aggregate state back into *state_data_type*. This function must take two arguments of types *bytea* and *internal*, and produce a result of type *internal*. (Note: the second, *internal* argument is unused, but is required for type safety reasons.)

initial_condition

The initial setting for the state value. This must be a string constant in the form accepted for the data type *state_data_type*. If not specified, the state value starts out null.

msfunc

The name of the forward state transition function to be called for each input row in moving-aggregate mode. This is exactly like the regular transition function, except that its first argument and result are of type *mstate_data_type*, which might be different from *state_data_type*.

minvfunc

The name of the inverse state transition function to be used in moving-aggregate mode. This function has the same argument and result types as *msfunc*, but it is used to remove a value from the current aggregate state, rather than add a value to it. The inverse transition function must have the same strictness attribute as the forward state transition function.

mstate_data_type

The data type for the aggregate's state value, when using moving-aggregate mode.

mstate_data_size

The approximate average size (in bytes) of the aggregate's state value, when using moving-aggregate mode. This works the same as *state_data_size*.

mffunc

The name of the final function called to compute the aggregate's result after all input rows have been traversed, when using moving-aggregate mode. This works the same as *ffunc*, except that

its first argument's type is `mstate_data_type` and extra dummy arguments are specified by writing `MFINALFUNC_EXTRA`. The aggregate result type determined by `mffunc` or `mstate_data_type` must match that determined by the aggregate's regular implementation.

`MFINALFUNC_MODIFY = { READ_ONLY | SHAREABLE | READ_WRITE }`

This option is like `FINALFUNC_MODIFY`, but it describes the behavior of the moving-aggregate final function.

minitial_condition

The initial setting for the state value, when using moving-aggregate mode. This works the same as *initial_condition*.

sort_operator

The associated sort operator for a `MIN`- or `MAX`-like aggregate. This is just an operator name (possibly schema-qualified). The operator is assumed to have the same input data types as the aggregate (which must be a single-argument normal aggregate).

`PARALLEL = { SAFE | RESTRICTED | UNSAFE }`

The meanings of `PARALLEL SAFE`, `PARALLEL RESTRICTED`, and `PARALLEL UNSAFE` are the same as in [CREATE FUNCTION](#). An aggregate will not be considered for parallelization if it is marked `PARALLEL UNSAFE` (which is the default!) or `PARALLEL RESTRICTED`. Note that the parallel-safety markings of the aggregate's support functions are not consulted by the planner, only the marking of the aggregate itself.

`HYPOTHETICAL`

For ordered-set aggregates only, this flag specifies that the aggregate arguments are to be processed according to the requirements for hypothetical-set aggregates: that is, the last few direct arguments must match the data types of the aggregated (`WITHIN GROUP`) arguments. The `HYPOTHETICAL` flag has no effect on run-time behavior, only on parse-time resolution of the data types and collations of the aggregate's arguments.

The parameters of `CREATE AGGREGATE` can be written in any order, not just the order illustrated above.

Notes

In parameters that specify support function names, you can write a schema name if needed, for example `SFUNC = public.sum`. Do not write argument types there, however — the argument types of the support functions are determined from other parameters.

Ordinarily, Postgres Pro functions are expected to be true functions that do not modify their input values. However, an aggregate transition function, *when used in the context of an aggregate*, is allowed to cheat and modify its transition-state argument in place. This can provide substantial performance benefits compared to making a fresh copy of the transition state each time.

Likewise, while an aggregate final function is normally expected not to modify its input values, sometimes it is impractical to avoid modifying the transition-state argument. Such behavior must be declared using the `FINALFUNC_MODIFY` parameter. The `READ_WRITE` value indicates that the final function modifies the transition state in unspecified ways. This value prevents use of the aggregate as a window function, and it also prevents merging of transition states for aggregate calls that share the same input values and transition functions. The `SHAREABLE` value indicates that the transition function cannot be applied after the final function, but multiple final-function calls can be performed on the ending transition state value. This value prevents use of the aggregate as a window function, but it allows merging of transition states. (That is, the optimization of interest here is not applying the same final function repeatedly, but applying different final functions to the same ending transition state value. This is allowed as long as none of the final functions are marked `READ_WRITE`.)

If an aggregate supports moving-aggregate mode, it will improve calculation efficiency when the aggregate is used as a window function for a window with moving frame start (that is, a frame start mode other

than `UNBOUNDED PRECEDING`). Conceptually, the forward transition function adds input values to the aggregate's state when they enter the window frame from the bottom, and the inverse transition function removes them again when they leave the frame at the top. So, when values are removed, they are always removed in the same order they were added. Whenever the inverse transition function is invoked, it will thus receive the earliest added but not yet removed argument value(s). The inverse transition function can assume that at least one row will remain in the current state after it removes the oldest row. (When this would not be the case, the window function mechanism simply starts a fresh aggregation, rather than using the inverse transition function.)

The forward transition function for moving-aggregate mode is not allowed to return `NULL` as the new state value. If the inverse transition function returns `NULL`, this is taken as an indication that the inverse function cannot reverse the state calculation for this particular input, and so the aggregate calculation will be redone from scratch for the current frame starting position. This convention allows moving-aggregate mode to be used in situations where there are some infrequent cases that are impractical to reverse out of the running state value.

If no moving-aggregate implementation is supplied, the aggregate can still be used with moving frames, but Postgres Pro will recompute the whole aggregation whenever the start of the frame moves. Note that whether or not the aggregate supports moving-aggregate mode, Postgres Pro can handle a moving frame end without recalculation; this is done by continuing to add new values to the aggregate's state. This is why use of an aggregate as a window function requires that the final function be read-only: it must not damage the aggregate's state value, so that the aggregation can be continued even after an aggregate result value has been obtained for one set of frame boundaries.

The syntax for ordered-set aggregates allows `VARIADIC` to be specified for both the last direct parameter and the last aggregated (`WITHIN GROUP`) parameter. However, the current implementation restricts use of `VARIADIC` in two ways. First, ordered-set aggregates can only use `VARIADIC "any"`, not other variadic array types. Second, if the last direct parameter is `VARIADIC "any"`, then there can be only one aggregated parameter and it must also be `VARIADIC "any"`. (In the representation used in the system catalogs, these two parameters are merged into a single `VARIADIC "any"` item, since `pg_proc` cannot represent functions with more than one `VARIADIC` parameter.) If the aggregate is a hypothetical-set aggregate, the direct arguments that match the `VARIADIC "any"` parameter are the hypothetical ones; any preceding parameters represent additional direct arguments that are not constrained to match the aggregated arguments.

Currently, ordered-set aggregates do not need to support moving-aggregate mode, since they cannot be used as window functions.

Partial (including parallel) aggregation is currently not supported for ordered-set aggregates. Also, it will never be used for aggregate calls that include `DISTINCT` or `ORDER BY` clauses, since those semantics cannot be supported during partial aggregation.

Examples

See [Section 39.12](#).

Compatibility

`CREATE AGGREGATE` is a Postgres Pro language extension. The SQL standard does not provide for user-defined aggregate functions.

See Also

[ALTER AGGREGATE](#), [DROP AGGREGATE](#)

CREATE CAST

CREATE CAST — define a new cast

Synopsis

```
CREATE CAST (source_type AS target_type)
  WITH FUNCTION function_name [ (argument_type [, ...]) ]
  [ AS ASSIGNMENT | AS IMPLICIT ]
```

```
CREATE CAST (source_type AS target_type)
  WITHOUT FUNCTION
  [ AS ASSIGNMENT | AS IMPLICIT ]
```

```
CREATE CAST (source_type AS target_type)
  WITH INOUT
  [ AS ASSIGNMENT | AS IMPLICIT ]
```

Description

CREATE CAST defines a new cast. A cast specifies how to perform a conversion between two data types. For example,

```
SELECT CAST(42 AS float8);
```

converts the integer constant 42 to type `float8` by invoking a previously specified function, in this case `float8(int4)`. (If no suitable cast has been defined, the conversion fails.)

Two types can be *binary coercible*, which means that the conversion can be performed “for free” without invoking any function. This requires that corresponding values use the same internal representation. For instance, the types `text` and `varchar` are binary coercible both ways. Binary coercibility is not necessarily a symmetric relationship. For example, the cast from `xml` to `text` can be performed for free in the present implementation, but the reverse direction requires a function that performs at least a syntax check. (Two types that are binary coercible both ways are also referred to as binary compatible.)

You can define a cast as an *I/O conversion cast* by using the `WITH INOUT` syntax. An I/O conversion cast is performed by invoking the output function of the source data type, and passing the resulting string to the input function of the target data type. In many common cases, this feature avoids the need to write a separate cast function for conversion. An I/O conversion cast acts the same as a regular function-based cast; only the implementation is different.

By default, a cast can be invoked only by an explicit cast request, that is an explicit `CAST(x AS typename)` or `x::typename` construct.

If the cast is marked `AS ASSIGNMENT` then it can be invoked implicitly when assigning a value to a column of the target data type. For example, supposing that `foo.f1` is a column of type `text`, then:

```
INSERT INTO foo (f1) VALUES (42);
```

will be allowed if the cast from type `integer` to type `text` is marked `AS ASSIGNMENT`, otherwise not. (We generally use the term *assignment cast* to describe this kind of cast.)

If the cast is marked `AS IMPLICIT` then it can be invoked implicitly in any context, whether assignment or internally in an expression. (We generally use the term *implicit cast* to describe this kind of cast.) For example, consider this query:

```
SELECT 2 + 4.0;
```

The parser initially marks the constants as being of type `integer` and `numeric` respectively. There is no `integer + numeric` operator in the system catalogs, but there is a `numeric + numeric` operator. The

query will therefore succeed if a cast from `integer` to `numeric` is available and is marked `AS IMPLICIT` — which in fact it is. The parser will apply the implicit cast and resolve the query as if it had been written

```
SELECT CAST ( 2 AS numeric ) + 4.0;
```

Now, the catalogs also provide a cast from `numeric` to `integer`. If that cast were marked `AS IMPLICIT` — which it is not — then the parser would be faced with choosing between the above interpretation and the alternative of casting the `numeric` constant to `integer` and applying the `integer + integer` operator. Lacking any knowledge of which choice to prefer, it would give up and declare the query ambiguous. The fact that only one of the two casts is implicit is the way in which we teach the parser to prefer resolution of a mixed `numeric-and-integer` expression as `numeric`; there is no built-in knowledge about that.

It is wise to be conservative about marking casts as implicit. An overabundance of implicit casting paths can cause Postgres Pro to choose surprising interpretations of commands, or to be unable to resolve commands at all because there are multiple possible interpretations. A good rule of thumb is to make a cast implicitly invocable only for information-preserving transformations between types in the same general type category. For example, the cast from `int2` to `int4` can reasonably be implicit, but the cast from `float8` to `int4` should probably be assignment-only. Cross-type-category casts, such as `text` to `int4`, are best made explicit-only.

Note

Sometimes it is necessary for usability or standards-compliance reasons to provide multiple implicit casts among a set of types, resulting in ambiguity that cannot be avoided as above. The parser has a fallback heuristic based on *type categories* and *preferred types* that can help to provide desired behavior in such cases. See [CREATE TYPE](#) for more information.

To be able to create a cast, you must own the source or the target data type and have `USAGE` privilege on the other type. To create a binary-coercible cast, you must be superuser. (This restriction is made because an erroneous binary-coercible cast conversion can easily crash the server.)

Parameters

source_type

The name of the source data type of the cast.

target_type

The name of the target data type of the cast.

function_name[(*argument_type* [, ...])]

The function used to perform the cast. The function name can be schema-qualified. If it is not, the function will be looked up in the schema search path. The function's result data type must match the target type of the cast. Its arguments are discussed below. If no argument list is specified, the function name must be unique in its schema.

`WITHOUT FUNCTION`

Indicates that the source type is binary-coercible to the target type, so no function is required to perform the cast.

`WITH INOUT`

Indicates that the cast is an I/O conversion cast, performed by invoking the output function of the source data type, and passing the resulting string to the input function of the target data type.

`AS ASSIGNMENT`

Indicates that the cast can be invoked implicitly in assignment contexts.

AS IMPLICIT

Indicates that the cast can be invoked implicitly in any context.

Cast implementation functions can have one to three arguments. The first argument type must be identical to or binary-coercible from the cast's source type. The second argument, if present, must be type `integer`; it receives the type modifier associated with the destination type, or `-1` if there is none. The third argument, if present, must be type `boolean`; it receives `true` if the cast is an explicit cast, `false` otherwise. (Bizarrely, the SQL standard demands different behaviors for explicit and implicit casts in some cases. This argument is supplied for functions that must implement such casts. It is not recommended that you design your own data types so that this matters.)

The return type of a cast function must be identical to or binary-coercible to the cast's target type.

Ordinarily a cast must have different source and target data types. However, it is allowed to declare a cast with identical source and target types if it has a cast implementation function with more than one argument. This is used to represent type-specific length coercion functions in the system catalogs. The named function is used to coerce a value of the type to the type modifier value given by its second argument.

When a cast has different source and target types and a function that takes more than one argument, it supports converting from one type to another and applying a length coercion in a single step. When no such entry is available, coercion to a type that uses a type modifier involves two cast steps, one to convert between data types and a second to apply the modifier.

A cast to or from a domain type currently has no effect. Casting to or from a domain uses the casts associated with its underlying type.

Notes

Use `DROP CAST` to remove user-defined casts.

Remember that if you want to be able to convert types both ways you need to declare casts both ways explicitly.

It is normally not necessary to create casts between user-defined types and the standard string types (`text`, `varchar`, and `char(n)`, as well as user-defined types that are defined to be in the string category). Postgres Pro provides automatic I/O conversion casts for that. The automatic casts to string types are treated as assignment casts, while the automatic casts from string types are explicit-only. You can override this behavior by declaring your own cast to replace an automatic cast, but usually the only reason to do so is if you want the conversion to be more easily invocable than the standard assignment-only or explicit-only setting. Another possible reason is that you want the conversion to behave differently from the type's I/O function; but that is sufficiently surprising that you should think twice about whether it's a good idea. (A small number of the built-in types do indeed have different behaviors for conversions, mostly because of requirements of the SQL standard.)

While not required, it is recommended that you continue to follow this old convention of naming cast implementation functions after the target data type. Many users are used to being able to cast data types using a function-style notation, that is `typename(x)`. This notation is in fact nothing more nor less than a call of the cast implementation function; it is not specially treated as a cast. If your conversion functions are not named to support this convention then you will have surprised users. Since Postgres Pro allows overloading of the same function name with different argument types, there is no difficulty in having multiple conversion functions from different types that all use the target type's name.

Note

Actually the preceding paragraph is an oversimplification: there are two cases in which a function-call construct will be treated as a cast request without having matched it to an actual function. If a function call `name(x)` does not exactly match any existing function, but `name` is the name of a

data type and `pg_cast` provides a binary-coercible cast to this type from the type of `x`, then the call will be construed as a binary-coercible cast. This exception is made so that binary-coercible casts can be invoked using functional syntax, even though they lack any function. Likewise, if there is no `pg_cast` entry but the cast would be to or from a string type, the call will be construed as an I/O conversion cast. This exception allows I/O conversion casts to be invoked using functional syntax.

Note

There is also an exception to the exception: I/O conversion casts from composite types to string types cannot be invoked using functional syntax, but must be written in explicit cast syntax (either `CAST` or `::` notation). This exception was added because after the introduction of automatically-provided I/O conversion casts, it was found too easy to accidentally invoke such a cast when a function or column reference was intended.

Examples

To create an assignment cast from type `bigint` to type `int4` using the function `int4(bigint)`:

```
CREATE CAST (bigint AS int4) WITH FUNCTION int4(bigint) AS ASSIGNMENT;
```

(This cast is already predefined in the system.)

Compatibility

The `CREATE CAST` command conforms to the SQL standard, except that SQL does not make provisions for binary-coercible types or extra arguments to implementation functions. `AS IMPLICIT` is a PostgreSQL extension, too.

See Also

[CREATE FUNCTION](#), [CREATE TYPE](#), [DROP CAST](#)

CREATE COLLATION

CREATE COLLATION — define a new collation

Synopsis

```
CREATE COLLATION [ IF NOT EXISTS ] name (  
    [ LOCALE = locale, ]  
    [ LC_COLLATE = lc_collate, ]  
    [ LC_CTYPE = lc_ctype, ]  
    [ PROVIDER = provider, ]  
    [ DETERMINISTIC = boolean, ]  
    [ RULES = rules, ]  
    [ VERSION = version ]  
)  
CREATE COLLATION [ IF NOT EXISTS ] name FROM existing_collation
```

Description

CREATE COLLATION defines a new collation using the specified operating system locale settings, or by copying an existing collation.

To be able to create a collation, you must have CREATE privilege on the destination schema.

Parameters

IF NOT EXISTS

Do not throw an error if a collation with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing collation is anything like the one that would have been created.

name

The name of the collation. The collation name can be schema-qualified. If it is not, the collation is defined in the current schema. The collation name must be unique within that schema. (The system catalogs can contain collations with the same name for other encodings, but these are ignored if the database encoding does not match.)

locale

The locale name for this collation. See [Section 23.2.2.3.1](#) and [Section 23.2.2.3.2](#) for details.

If *provider* is `libc`, this is a shortcut for setting `LC_COLLATE` and `LC_CTYPE` at once. If you specify *locale*, you cannot specify either of those parameters.

If *provider* is `builtin`, then *locale* must be specified and set to either `C` or `C.UTF-8`.

lc_collate

If *provider* is `libc`, use the specified operating system locale for the `LC_COLLATE` locale category.

lc_ctype

If *provider* is `libc`, use the specified operating system locale for the `LC_CTYPE` locale category.

provider

Specifies the provider to use for locale services associated with this collation. Possible values are `builtin`, `icu` (if the server was built with ICU support) or `libc`. `libc` is the default. See [Section 23.1.4](#) for details.

DETERMINISTIC

Specifies whether the collation should use deterministic comparisons. The default is true. A deterministic comparison considers strings that are not byte-wise equal to be unequal even if they are considered logically equal by the comparison. Postgres Pro breaks ties using a byte-wise comparison. Comparison that is not deterministic can make the collation be, say, case- or accent-insensitive. For that, you need to choose an appropriate `LOCALE` setting *and* set the collation to not deterministic here.

Nondeterministic collations are only supported with the ICU provider.

rules

Specifies additional collation rules to customize the behavior of the collation. This is supported for ICU only. See [Section 23.2.3.4](#) for details.

version

Specifies the version string to store with the collation. Normally, this should be omitted, which will cause the version to be computed from the actual version of the collation as provided by the operating system. This option is intended to be used by `pg_upgrade` for copying the version from an existing installation.

See also [ALTER COLLATION](#) for how to handle collation version mismatches.

existing_collation

The name of an existing collation to copy. The new collation will have the same properties as the existing one, but it will be an independent object.

Notes

`CREATE COLLATION` takes a `SHARE ROW EXCLUSIVE` lock, which is self-conflicting, on the `pg_collation` system catalog, so only one `CREATE COLLATION` command can run at a time.

Use `DROP COLLATION` to remove user-defined collations.

See [Section 23.2.2.3](#) for more information on how to create collations.

When using the `libc` collation provider, the locale must be applicable to the current database encoding. See [CREATE DATABASE](#) for the precise rules.

Examples

To create a collation from the operating system locale `ru_RU.utf8` (assuming the current database encoding is UTF8):

```
CREATE COLLATION russian (locale = 'ru_RU.utf8');
```

To create a collation using the ICU provider where Latin characters precede Cyrillic ones:

```
CREATE COLLATION latn_cyrl (provider = icu, locale = 'ru-RU-u-kr-latn-cyrl');
```

To create a collation using the ICU provider, based on the root ICU locale, with custom rules:

```
CREATE COLLATION custom (provider = icu, locale = 'und', rules = '&V << w <<< W');
```

See [Section 23.2.3.4](#) for further details and examples on the rules syntax.

To create a collation from an existing collation:

```
CREATE COLLATION german FROM "de_DE";
```

This can be convenient to be able to use operating-system-independent collation names in applications.

Compatibility

There is a `CREATE COLLATION` statement in the SQL standard, but it is limited to copying an existing collation. The syntax to create a new collation is a Postgres Pro extension.

See Also

[ALTER COLLATION](#), [DROP COLLATION](#)

CREATE CONVERSION

CREATE CONVERSION — define a new encoding conversion

Synopsis

```
CREATE [ DEFAULT ] CONVERSION name
    FOR source_encoding TO dest_encoding FROM function_name
```

Description

CREATE CONVERSION defines a new conversion between two character set encodings.

Conversions that are marked `DEFAULT` can be used for automatic encoding conversion between client and server. To support that usage, two conversions, from encoding A to B *and* from encoding B to A, must be defined.

To be able to create a conversion, you must have `EXECUTE` privilege on the function and `CREATE` privilege on the destination schema.

Parameters

DEFAULT

The `DEFAULT` clause indicates that this conversion is the default for this particular source to destination encoding. There should be only one default encoding in a schema for the encoding pair.

name

The name of the conversion. The conversion name can be schema-qualified. If it is not, the conversion is defined in the current schema. The conversion name must be unique within a schema.

source_encoding

The source encoding name.

dest_encoding

The destination encoding name.

function_name

The function used to perform the conversion. The function name can be schema-qualified. If it is not, the function will be looked up in the path.

The function must have the following signature:

```
conv_proc (
    integer,  -- source encoding ID
    integer,  -- destination encoding ID
    cstring,  -- source string (null terminated C string)
    internal, -- destination (fill with a null terminated C string)
    integer,  -- source string length
    boolean   -- if true, don't throw an error if conversion fails
) RETURNS integer;
```

The return value is the number of source bytes that were successfully converted. If the last argument is false, the function must throw an error on invalid input, and the return value is always equal to the source string length.

Notes

Neither the source nor the destination encoding can be `SQL_ASCII`, as the server's behavior for cases involving the `SQL_ASCII` “encoding” is hard-wired.

Use `DROP CONVERSION` to remove user-defined conversions.

The privileges required to create a conversion might be changed in a future release.

Examples

To create a conversion from encoding `UTF8` to `LATIN1` using `myfunc`:

```
CREATE CONVERSION myconv FOR 'UTF8' TO 'LATIN1' FROM myfunc;
```

Compatibility

`CREATE CONVERSION` is a Postgres Pro extension. There is no `CREATE CONVERSION` statement in the SQL standard, but a `CREATE TRANSLATION` statement that is very similar in purpose and syntax.

See Also

[ALTER CONVERSION](#), [CREATE FUNCTION](#), [DROP CONVERSION](#)

CREATE DATABASE

CREATE DATABASE — create a new database

Synopsis

```
CREATE DATABASE name
    [ WITH ] [ OWNER [=] user_name ]
    [ TEMPLATE [=] template ]
    [ ENCODING [=] encoding ]
    [ STRATEGY [=] strategy ]
    [ LOCALE [=] locale ]
    [ LC_COLLATE [=] lc_collate ]
    [ LC_CTYPE [=] lc_ctype ]
    [ BUILTIN_LOCALE [=] builtin_locale ]
    [ ICU_LOCALE [=] icu_locale ]
    [ ICU_RULES [=] icu_rules ]
    [ LOCALE_PROVIDER [=] locale_provider ]
    [ COLLATION_VERSION = collation_version ]
    [ TABLESPACE [=] tablespace_name ]
    [ ALLOW_CONNECTIONS [=] allowconn ]
    [ CONNECTION LIMIT [=] conlimit ]
    [ IS_TEMPLATE [=] istemplate ]
    [ OID [=] oid ]
```

Description

CREATE DATABASE creates a new Postgres Pro database.

To create a database, you must be a superuser or have the special `CREATEDB` privilege. See [CREATE ROLE](#).

By default, the new database will be created by cloning the standard system database `template1`. A different template can be specified by writing `TEMPLATE name`. In particular, by writing `TEMPLATE template0`, you can create a pristine database (one where no user-defined objects exist and where the system objects have not been altered) containing only the standard objects predefined by your version of Postgres Pro. This is useful if you wish to avoid copying any installation-local objects that might have been added to `template1`.

Parameters

name

The name of a database to create.

user_name

The role name of the user who will own the new database, or `DEFAULT` to use the default (namely, the user executing the command). To create a database owned by another role, you must be able to `SET ROLE` to that role.

template

The name of the template from which to create the new database, or `DEFAULT` to use the default template (`template1`).

encoding

Character set encoding to use in the new database. Specify a string constant (e.g., `'SQL_ASCII'`), or an integer encoding number, or `DEFAULT` to use the default encoding (namely, the encoding of

the template database). The character sets supported by the Postgres Pro server are described in [Section 23.3.1](#). See below for additional restrictions.

strategy

Strategy to be used in creating the new database. If the `WAL_LOG` strategy is used, the database will be copied block by block and each block will be separately written to the write-ahead log. This is the most efficient strategy in cases where the template database is small, and therefore it is the default. The older `FILE_COPY` strategy is also available. This strategy writes a small record to the write-ahead log for each tablespace used by the target database. Each such record represents copying an entire directory to a new location at the filesystem level. While this does reduce the write-ahead log volume substantially, especially if the template database is large, it also forces the system to perform a checkpoint both before and after the creation of the new database. In some situations, this may have a noticeable negative impact on overall system performance.

locale

Sets the default collation order and character classification in the new database. Collation affects the sort order applied to strings, e.g., in queries with `ORDER BY`, as well as the order used in indexes on text columns. Character classification affects the categorization of characters, e.g., lower, upper, and digit. Also sets the associated aspects of the operating system environment, `LC_COLLATE` and `LC_CTYPE`. The default is the same setting as the template database. See [Section 23.2.2.3.1](#) and [Section 23.2.2.3.2](#) for details.

Can be overridden by setting `lc_collate`, `lc_ctype`, `builtin_locale`, or `icu_locale` individually.

If `locale_provider` is `builtin`, then `locale` or `builtin_locale` must be specified and set to either `C` or `C.UTF-8`.

Tip

The other locale settings `lc_messages`, `lc_monetary`, `lc_numeric`, and `lc_time` are not fixed per database and are not set by this command. If you want to make them the default for a specific database, you can use `ALTER DATABASE ... SET`.

lc_collate

Sets `LC_COLLATE` in the database server's operating system environment. The default is the setting of `locale` if specified, otherwise the same setting as the template database. See below for additional restrictions.

If `locale_provider` is `libc`, also sets the default collation order to use in the new database, overriding the setting `locale`.

lc_ctype

Sets `LC_CTYPE` in the database server's operating system environment. The default is the setting of `locale` if specified, otherwise the same setting as the template database. See below for additional restrictions.

If `locale_provider` is `libc`, also sets the default character classification to use in the new database, overriding the setting `locale`.

builtin_locale

Specifies the builtin provider locale for the database default collation order and character classification, overriding the setting `locale`. The `locale provider` must be `builtin`. The default is the setting of `locale` if specified; otherwise the same setting as the template database.

The locales available for the `builtin` provider are `C` and `C.UTF-8`.

icu_locale

Specifies the ICU locale (see [Section 23.2.2.3.2](#)) for the database default collation order and character classification, overriding the setting *locale*. The *locale provider* must be ICU. The default is the setting of *locale* if specified; otherwise the same setting as the template database.

icu_rules

Specifies additional collation rules to customize the behavior of the default collation of this database. This is supported for ICU only. See [Section 23.2.3.4](#) for details.

locale_provider

Specifies the provider to use for the default collation in this database. Possible values are *builtin*, *icu* (if the server was built with ICU support) or *libc*. By default, the provider is the same as that of the *template*. See [Section 23.1.4](#) for details.

collation_version

Specifies the collation version string to store with the database. Normally, this should be omitted, which will cause the version to be computed from the actual version of the database collation as provided by the operating system. This option is intended to be used by *pg_upgrade* for copying the version from an existing installation.

See also [ALTER DATABASE](#) for how to handle database collation version mismatches.

tablespace_name

The name of the tablespace that will be associated with the new database, or *DEFAULT* to use the template database's tablespace. This tablespace will be the default tablespace used for objects created in this database. See [CREATE TABLESPACE](#) for more information.

allowconn

If false then no one can connect to this database. The default is true, allowing connections (except as restricted by other mechanisms, such as *GRANT/REVOKE CONNECT*).

connlimit

How many concurrent connections can be made to this database. -1 (the default) means no limit.

istemplate

If true, then this database can be cloned by any user with *CREATEDB* privileges; if false (the default), then only superusers or the owner of the database can clone it.

oid

The object identifier to be used for the new database. If this parameter is not specified, Postgres Pro will choose a suitable OID automatically. This parameter is primarily intended for internal use by *pg_upgrade*, and only *pg_upgrade* can specify a value less than 16384.

Optional parameters can be written in any order, not only the order illustrated above.

Notes

CREATE DATABASE cannot be executed inside a transaction block.

Errors along the line of “could not initialize database directory” are most likely related to insufficient permissions on the data directory, a full disk, or other file system problems.

Use [DROP DATABASE](#) to remove a database.

The program `createdb` is a wrapper program around this command, provided for convenience.

Database-level configuration parameters (set via `ALTER DATABASE`) and database-level permissions (set via `GRANT`) are not copied from the template database.

Although it is possible to copy a database other than `template1` by specifying its name as the template, this is not (yet) intended as a general-purpose “COPY DATABASE” facility. The principal limitation is that no other sessions can be connected to the template database while it is being copied. `CREATE DATABASE` will fail if any other connection exists when it starts; otherwise, new connections to the template database are locked out until `CREATE DATABASE` completes. See [Section 22.4](#) for more information.

The character set encoding specified for the new database must be compatible with the chosen locale settings (`LC_COLLATE` and `LC_CTYPE`). If the locale is `C` (or equivalently `POSIX`), then all encodings are allowed, but for other locale settings there is only one encoding that will work properly. (On Windows, however, UTF-8 encoding can be used with any locale.) `CREATE DATABASE` will allow superusers to specify `SQL_ASCII` encoding regardless of the locale settings, but this choice is deprecated and may result in misbehavior of character-string functions if data that is not encoding-compatible with the locale is stored in the database.

The encoding and locale settings must match those of the template database, except when `template0` is used as template. This is because other databases might contain data that does not match the specified encoding, or might contain indexes whose sort ordering is affected by `LC_COLLATE` and `LC_CTYPE`. Copying such data would result in a database that is corrupt according to the new settings. `template0`, however, is known to not contain any data or indexes that would be affected.

There is currently no option to use a database locale with nondeterministic comparisons (see [CREATE COLLATION](#) for an explanation). If this is needed, then per-column collations would need to be used.

The `CONNECTION LIMIT` option is only enforced approximately; if two new sessions start at about the same time when just one connection “slot” remains for the database, it is possible that both will fail. Also, the limit is not enforced against superusers or background worker processes.

Examples

To create a new database:

```
CREATE DATABASE lusiadas;
```

To create a database `sales` owned by user `salesapp` with a default tablespace of `salesspace`:

```
CREATE DATABASE sales OWNER salesapp TABLESPACE salesspace;
```

To create a database `music` with a different locale:

```
CREATE DATABASE music
    LOCALE 'sv_SE.utf8'
    TEMPLATE template0;
```

In this example, the `TEMPLATE template0` clause is required if the specified locale is different from the one in `template1`. (If it is not, then specifying the locale explicitly is redundant.)

To create a database `music2` with a different locale and a different character set encoding:

```
CREATE DATABASE music2
    LOCALE 'sv_SE.iso885915'
    ENCODING LATIN9
    TEMPLATE template0;
```

The specified locale and encoding settings must match, or an error will be reported.

Note that locale names are specific to the operating system, so that the above commands might not work in the same way everywhere.

Compatibility

There is no `CREATE DATABASE` statement in the SQL standard. Databases are equivalent to catalogs, whose creation is implementation-defined.

See Also

[ALTER DATABASE](#), [DROP DATABASE](#)

CREATE DOMAIN

CREATE DOMAIN — define a new domain

Synopsis

```
CREATE DOMAIN name [ AS ] data_type
    [ COLLATE collation ]
    [ DEFAULT expression ]
    [ domain_constraint [ ... ] ]
```

where *domain_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ NOT NULL | NULL | CHECK (expression) }
```

Description

CREATE DOMAIN creates a new domain. A domain is essentially a data type with optional constraints (restrictions on the allowed set of values). The user who defines a domain becomes its owner.

If a schema name is given (for example, CREATE DOMAIN myschema.mydomain ...) then the domain is created in the specified schema. Otherwise it is created in the current schema. The domain name must be unique among the types and domains existing in its schema.

Domains are useful for abstracting common constraints on fields into a single location for maintenance. For example, several tables might contain email address columns, all requiring the same CHECK constraint to verify the address syntax. Define a domain rather than setting up each table's constraint individually.

To be able to create a domain, you must have USAGE privilege on the underlying type.

Parameters

name

The name (optionally schema-qualified) of a domain to be created.

data_type

The underlying data type of the domain. This can include array specifiers.

collation

An optional collation for the domain. If no collation is specified, the domain has the same collation behavior as its underlying data type. The underlying type must be collatable if COLLATE is specified.

DEFAULT *expression*

The DEFAULT clause specifies a default value for columns of the domain data type. The value is any variable-free expression (but subqueries are not allowed). The data type of the default expression must match the data type of the domain. If no default value is specified, then the default value is the null value.

The default expression will be used in any insert operation that does not specify a value for the column. If a default value is defined for a particular column, it overrides any default associated with the domain. In turn, the domain default overrides any default value associated with the underlying data type.

CONSTRAINT *constraint_name*

An optional name for a constraint. If not specified, the system generates a name.

`NOT NULL`

Values of this domain are prevented from being null (but see notes below).

`NULL`

Values of this domain are allowed to be null. This is the default.

This clause is only intended for compatibility with nonstandard SQL databases. Its use is discouraged in new applications.

`CHECK (expression)`

`CHECK` clauses specify integrity constraints or tests which values of the domain must satisfy. Each constraint must be an expression producing a Boolean result. It should use the key word `VALUE` to refer to the value being tested. Expressions evaluating to `TRUE` or `UNKNOWN` succeed. If the expression produces a `FALSE` result, an error is reported and the value is not allowed to be converted to the domain type.

Currently, `CHECK` expressions cannot contain subqueries nor refer to variables other than `VALUE`.

When a domain has multiple `CHECK` constraints, they will be tested in alphabetical order by name. (PostgreSQL versions before 9.5 did not honor any particular firing order for `CHECK` constraints.)

Notes

Domain constraints, particularly `NOT NULL`, are checked when converting a value to the domain type. It is possible for a column that is nominally of the domain type to read as null despite there being such a constraint. For example, this can happen in an outer-join query, if the domain column is on the nullable side of the outer join. A more subtle example is

```
INSERT INTO tab (domcol) VALUES ((SELECT domcol FROM tab WHERE false));
```

The empty scalar sub-`SELECT` will produce a null value that is considered to be of the domain type, so no further constraint checking is applied to it, and the insertion will succeed.

It is very difficult to avoid such problems, because of SQL's general assumption that a null value is a valid value of every data type. Best practice therefore is to design a domain's constraints so that a null value is allowed, and then to apply column `NOT NULL` constraints to columns of the domain type as needed, rather than directly to the domain type.

Postgres Pro assumes that `CHECK` constraints' conditions are immutable, that is, they will always give the same result for the same input value. This assumption is what justifies examining `CHECK` constraints only when a value is first converted to be of a domain type, and not at other times. (This is essentially the same as the treatment of table `CHECK` constraints, as described in [Section 5.5.1](#).)

An example of a common way to break this assumption is to reference a user-defined function in a `CHECK` expression, and then change the behavior of that function. Postgres Pro does not disallow that, but it will not notice if there are stored values of the domain type that now violate the `CHECK` constraint. That would cause a subsequent database dump and restore to fail. The recommended way to handle such a change is to drop the constraint (using `ALTER DOMAIN`), adjust the function definition, and re-add the constraint, thereby rechecking it against stored data.

It's also good practice to ensure that domain `CHECK` expressions will not throw errors.

Examples

This example creates the `us_postal_code` data type and then uses the type in a table definition. A regular expression test is used to verify that the value looks like a valid US postal code:

```
CREATE DOMAIN us_postal_code AS TEXT
CHECK (
    VALUE ~ '^\d{5}$'
```

```
OR VALUE ~ '^\\d{5}-\\d{4}$'
);
```

```
CREATE TABLE us_snail_addy (
  address_id SERIAL PRIMARY KEY,
  street1 TEXT NOT NULL,
  street2 TEXT,
  street3 TEXT,
  city TEXT NOT NULL,
  postal us_postal_code NOT NULL
);
```

Compatibility

The command `CREATE DOMAIN` conforms to the SQL standard.

The syntax `NOT NULL` in this command is a Postgres Pro extension. (A standard-conforming way to write the same for non-composite data types would be `CHECK (VALUE IS NOT NULL)`. However, per [the section called “Notes”](#), such constraints are best avoided in practice anyway.) The `NULL` “constraint” is a Postgres Pro extension (see also [Compatibility](#)).

See Also

[ALTER DOMAIN](#), [DROP DOMAIN](#)

CREATE EVENT TRIGGER

CREATE EVENT TRIGGER — define a new event trigger

Synopsis

```
CREATE EVENT TRIGGER name
  ON event
  [ WHEN filter_variable IN (filter_value [, ... ]) [ AND ... ] ]
  EXECUTE { FUNCTION | PROCEDURE } function_name()
```

Description

CREATE EVENT TRIGGER creates a new event trigger. Whenever the designated event occurs and the WHEN condition associated with the trigger, if any, is satisfied, the trigger function will be executed. For a general introduction to event triggers, see [Chapter 41](#). The user who creates an event trigger becomes its owner.

Parameters

name

The name to give the new trigger. This name must be unique within the database.

event

The name of the event that triggers a call to the given function. See [Section 41.1](#) for more information on event names.

filter_variable

The name of a variable used to filter events. This makes it possible to restrict the firing of the trigger to a subset of the cases in which it is supported. Currently the only supported *filter_variable* is TAG.

filter_value

A list of values for the associated *filter_variable* for which the trigger should fire. For TAG, this means a list of command tags (e.g., 'DROP FUNCTION').

function_name

A user-supplied function that is declared as taking no argument and returning type event_trigger.

In the syntax of CREATE EVENT TRIGGER, the keywords FUNCTION and PROCEDURE are equivalent, but the referenced function must in any case be a function, not a procedure. The use of the keyword PROCEDURE here is historical and deprecated.

Notes

Only superusers can create event triggers.

Event triggers are disabled in single-user mode (see [postgres](#)) as well as when [event_triggers](#) is set to false. If an erroneous event trigger disables the database so much that you can't even drop the trigger, restart with [event_triggers](#) set to false to temporarily disable event triggers, or in single-user mode, and you'll be able to do that.

Examples

Forbid the execution of any [DDL](#) command:


```
CREATE OR REPLACE FUNCTION abort_any_command()
  RETURNS event_trigger
  LANGUAGE plpgsql
  AS $$
BEGIN
  RAISE EXCEPTION 'command % is disabled', tg_tag;
END;
$$;

CREATE EVENT TRIGGER abort_ddl ON ddl_command_start
  EXECUTE FUNCTION abort_any_command();
```

Compatibility

There is no `CREATE EVENT TRIGGER` statement in the SQL standard.

See Also

[ALTER EVENT TRIGGER](#), [DROP EVENT TRIGGER](#), [CREATE FUNCTION](#)

CREATE EXTENSION

CREATE EXTENSION — install an extension

Synopsis

```
CREATE EXTENSION [ IF NOT EXISTS ] extension_name
    [ WITH ] [ SCHEMA schema_name ]
    [ VERSION version ]
    [ CASCADE ]
```

Description

`CREATE EXTENSION` loads a new extension into the current database. There must not be an extension of the same name already loaded.

Loading an extension essentially amounts to running the extension's script file. The script will typically create new SQL objects such as functions, data types, operators and index support methods. `CREATE EXTENSION` additionally records the identities of all the created objects, so that they can be dropped again if `DROP EXTENSION` is issued.

The user who runs `CREATE EXTENSION` becomes the owner of the extension for purposes of later privilege checks, and normally also becomes the owner of any objects created by the extension's script.

Loading an extension ordinarily requires the same privileges that would be required to create its component objects. For many extensions this means superuser privileges are needed. However, if the extension is marked *trusted* in its control file, then it can be installed by any user who has `CREATE` privilege on the current database. In this case the extension object itself will be owned by the calling user, but the contained objects will be owned by the bootstrap superuser (unless the extension's script explicitly assigns them to the calling user). This configuration gives the calling user the right to drop the extension, but not to modify individual objects within it.

Parameters

`IF NOT EXISTS`

Do not throw an error if an extension with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing extension is anything like the one that would have been created from the currently-available script file.

extension_name

The name of the extension to be installed. Postgres Pro will create the extension using details from the file `SHAREDIR/extension/extension_name.control`.

schema_name

The name of the schema in which to install the extension's objects, given that the extension allows its contents to be relocated. The named schema must already exist. If not specified, and the extension's control file does not specify a schema either, the current default object creation schema is used.

If the extension specifies a `schema` parameter in its control file, then that schema cannot be overridden with a `SCHEMA` clause. Normally, an error will be raised if a `SCHEMA` clause is given and it conflicts with the extension's `schema` parameter. However, if the `CASCADE` clause is also given, then *schema_name* is ignored when it conflicts. The given *schema_name* will be used for installation of any needed extensions that do not specify `schema` in their control files.

Remember that the extension itself is not considered to be within any schema: extensions have unqualified names that must be unique database-wide. But objects belonging to the extension can be within schemas.

version

The version of the extension to install. This can be written as either an identifier or a string literal. The default version is whatever is specified in the extension's control file.

CASCADE

Automatically install any extensions that this extension depends on that are not already installed. Their dependencies are likewise automatically installed, recursively. The `SCHEMA` clause, if given, applies to all extensions that get installed this way. Other options of the statement are not applied to automatically-installed extensions; in particular, their default versions are always selected.

Notes

Before you can use `CREATE EXTENSION` to load an extension into a database, the extension's supporting files must be installed. Information about installing the extensions supplied with Postgres Pro can be found in [Additional Supplied Modules](#).

The extensions currently available for loading can be identified from the `pg_available_extensions` or `pg_available_extension_versions` system views.

Caution

Installing an extension as superuser requires trusting that the extension's author wrote the extension installation script in a secure fashion. It is not terribly difficult for a malicious user to create trojan-horse objects that will compromise later execution of a carelessly-written extension script, allowing that user to acquire superuser privileges. However, trojan-horse objects are only hazardous if they are in the `search_path` during script execution, meaning that they are in the extension's installation target schema or in the schema of some extension it depends on. Therefore, a good rule of thumb when dealing with extensions whose scripts have not been carefully vetted is to install them only into schemas for which `CREATE` privilege has not been and will not be granted to any untrusted users. Likewise for any extensions they depend on.

The extensions supplied with Postgres Pro are believed to be secure against installation-time attacks of this sort, except for a few that depend on other extensions. As stated in the documentation for those extensions, they should be installed into secure schemas, or installed into the same schemas as the extensions they depend on, or both.

For information about writing new extensions, see [Section 39.17](#).

Compatibility

`CREATE EXTENSION` is a Postgres Pro extension.

See Also

[ALTER EXTENSION](#), [DROP EXTENSION](#)

CREATE FUNCTION

CREATE FUNCTION — define a new function

Synopsis

```
CREATE [ OR REPLACE ] FUNCTION
    name ( [ [ argmode ] [ argname ] argtype [ { DEFAULT | = } default_expr ]
    [, ...] ] )
    [ RETURNS rettype
      | RETURNS TABLE ( column_name column_type [, ...] ) ]
    { LANGUAGE lang_name
      | TRANSFORM { FOR TYPE type_name } [, ... ]
      | WINDOW
      | { IMMUTABLE | STABLE | VOLATILE }
      | [ NOT ] LEAKPROOF
      | { CALLED ON NULL INPUT | RETURNS NULL ON NULL INPUT | STRICT }
      | { [ EXTERNAL ] SECURITY INVOKER | [ EXTERNAL ] SECURITY DEFINER }
      | PARALLEL { UNSAFE | RESTRICTED | SAFE }
      | COST execution_cost
      | ROWS result_rows
      | SUPPORT support_function
      | SET configuration_parameter { TO value | = value | FROM CURRENT }
      | AS 'definition'
      | AS 'obj_file', 'link_symbol'
      | sql_body
    } ...
```

Description

`CREATE FUNCTION` defines a new function. `CREATE OR REPLACE FUNCTION` will either create a new function, or replace an existing definition. To be able to define a function, the user must have the `USAGE` privilege on the language.

If a schema name is included, then the function is created in the specified schema. Otherwise it is created in the current schema. The name of the new function must not match any existing function or procedure with the same input argument types in the same schema. However, functions and procedures of different argument types can share a name (this is called *overloading*).

To replace the current definition of an existing function, use `CREATE OR REPLACE FUNCTION`. It is not possible to change the name or argument types of a function this way (if you tried, you would actually be creating a new, distinct function). Also, `CREATE OR REPLACE FUNCTION` will not let you change the return type of an existing function. To do that, you must drop and recreate the function. (When using `OUT` parameters, that means you cannot change the types of any `OUT` parameters except by dropping the function.)

When `CREATE OR REPLACE FUNCTION` is used to replace an existing function, the ownership and permissions of the function do not change. All other function properties are assigned the values specified or implied in the command. You must own the function to replace it (this includes being a member of the owning role).

If you drop and then recreate a function, the new function is not the same entity as the old; you will have to drop existing rules, views, triggers, etc. that refer to the old function. Use `CREATE OR REPLACE FUNCTION` to change a function definition without breaking objects that refer to the function. Also, `ALTER FUNCTION` can be used to change most of the auxiliary properties of an existing function.

The user that creates the function becomes the owner of the function.

To be able to create a function, you must have `USAGE` privilege on the argument types and the return type.

Refer to [Section 39.3](#) for further information on writing functions.

Parameters

name

The name (optionally schema-qualified) of the function to create.

argmode

The mode of an argument: `IN`, `OUT`, `INOUT`, or `VARIADIC`. If omitted, the default is `IN`. Only `OUT` arguments can follow a `VARIADIC` one. Also, `OUT` and `INOUT` arguments cannot be used together with the `RETURNS TABLE` notation.

argname

The name of an argument. Some languages (including SQL and PL/pgSQL) let you use the name in the function body. For other languages the name of an input argument is just extra documentation, so far as the function itself is concerned; but you can use input argument names when calling a function to improve readability (see [Section 4.3](#)). In any case, the name of an output argument is significant, because it defines the column name in the result row type. (If you omit the name for an output argument, the system will choose a default column name.)

argtype

The data type(s) of the function's arguments (optionally schema-qualified), if any. The argument types can be base, composite, or domain types, or can reference the type of a table column.

Depending on the implementation language it might also be allowed to specify “pseudo-types” such as `cstring`. Pseudo-types indicate that the actual argument type is either incompletely specified, or outside the set of ordinary SQL data types.

The type of a column is referenced by writing `table_name.column_name%TYPE`. Using this feature can sometimes help make a function independent of changes to the definition of a table.

default_expr

An expression to be used as default value if the parameter is not specified. The expression has to be coercible to the argument type of the parameter. Only input (including `INOUT`) parameters can have a default value. All input parameters following a parameter with a default value must have default values as well.

rettype

The return data type (optionally schema-qualified). The return type can be a base, composite, or domain type, or can reference the type of a table column. Depending on the implementation language it might also be allowed to specify “pseudo-types” such as `cstring`. If the function is not supposed to return a value, specify `void` as the return type.

When there are `OUT` or `INOUT` parameters, the `RETURNS` clause can be omitted. If present, it must agree with the result type implied by the output parameters: `RECORD` if there are multiple output parameters, or the same type as the single output parameter.

The `SETOF` modifier indicates that the function will return a set of items, rather than a single item.

The type of a column is referenced by writing `table_name.column_name%TYPE`.

column_name

The name of an output column in the `RETURNS TABLE` syntax. This is effectively another way of declaring a named `OUT` parameter, except that `RETURNS TABLE` also implies `RETURNS SETOF`.

column_type

The data type of an output column in the `RETURNS TABLE` syntax.

lang_name

The name of the language that the function is implemented in. It can be `sql`, `c`, `internal`, or the name of a user-defined procedural language, e.g., `plpgsql`. The default is `sql` if *sql_body* is specified. Enclosing the name in single quotes is deprecated and requires matching case.

`TRANSFORM { FOR TYPE type_name } [, ...] }`

Lists which transforms a call to the function should apply. Transforms convert between SQL types and language-specific data types; see [CREATE TRANSFORM](#). Procedural language implementations usually have hardcoded knowledge of the built-in types, so those don't need to be listed here. If a procedural language implementation does not know how to handle a type and no transform is supplied, it will fall back to a default behavior for converting data types, but this depends on the implementation.

`WINDOW`

`WINDOW` indicates that the function is a *window function* rather than a plain function. This is currently only useful for functions written in C. The `WINDOW` attribute cannot be changed when replacing an existing function definition.

`IMMUTABLE`

`STABLE`

`VOLATILE`

These attributes inform the query optimizer about the behavior of the function. At most one choice can be specified. If none of these appear, `VOLATILE` is the default assumption.

`IMMUTABLE` indicates that the function cannot modify the database and always returns the same result when given the same argument values; that is, it does not do database lookups or otherwise use information not directly present in its argument list. If this option is given, any call of the function with all-constant arguments can be immediately replaced with the function value.

`STABLE` indicates that the function cannot modify the database, and that within a single table scan it will consistently return the same result for the same argument values, but that its result could change across SQL statements. This is the appropriate selection for functions whose results depend on database lookups, parameter variables (such as the current time zone), etc. (It is inappropriate for `AFTER` triggers that wish to query rows modified by the current command.) Also note that the `current_timestamp` family of functions qualify as stable, since their values do not change within a transaction.

`VOLATILE` indicates that the function value can change even within a single table scan, so no optimizations can be made. Relatively few database functions are volatile in this sense; some examples are `random()`, `currval()`, `timeofday()`. But note that any function that has side-effects must be classified volatile, even if its result is quite predictable, to prevent calls from being optimized away; an example is `setval()`.

For additional details see [Section 39.7](#).

`LEAKPROOF`

`LEAKPROOF` indicates that the function has no side effects. It reveals no information about its arguments other than by its return value. For example, a function which throws an error message for some argument values but not others, or which includes the argument values in any error message, is not leakproof. This affects how the system executes queries against views created with the `security_barrier` option or tables with row level security enabled. The system will enforce conditions from security policies and security barrier views before any user-supplied conditions from the query itself that contain non-leakproof functions, in order to prevent the inadvertent exposure of data.

Functions and operators marked as leakproof are assumed to be trustworthy, and may be executed before conditions from security policies and security barrier views. In addition, functions which do not take arguments or which are not passed any arguments from the security barrier view or table do not have to be marked as leakproof to be executed before security conditions. See [CREATE VIEW](#) and [Section 42.5](#). This option can only be set by the superuser.

`CALLED ON NULL INPUT`

`RETURNS NULL ON NULL INPUT`

`STRICT`

`CALLED ON NULL INPUT` (the default) indicates that the function will be called normally when some of its arguments are null. It is then the function author's responsibility to check for null values if necessary and respond appropriately.

`RETURNS NULL ON NULL INPUT` or `STRICT` indicates that the function always returns null whenever any of its arguments are null. If this parameter is specified, the function is not executed when there are null arguments; instead a null result is assumed automatically.

`[EXTERNAL] SECURITY INVOKER`

`[EXTERNAL] SECURITY DEFINER`

`SECURITY INVOKER` indicates that the function is to be executed with the privileges of the user that calls it. That is the default. `SECURITY DEFINER` specifies that the function is to be executed with the privileges of the user that owns it. For information on how to write `SECURITY DEFINER` functions safely, [see below](#).

The key word `EXTERNAL` is allowed for SQL conformance, but it is optional since, unlike in SQL, this feature applies to all functions not only external ones.

`PARALLEL`

`PARALLEL UNSAFE` indicates that the function can't be executed in parallel mode; the presence of such a function in an SQL statement forces a serial execution plan. This is the default. `PARALLEL RESTRICTED` indicates that the function can be executed in parallel mode, but only in the parallel group leader process. `PARALLEL SAFE` indicates that the function is safe to run in parallel mode without restriction, including in parallel worker processes.

Functions should be labeled parallel unsafe if they modify any database state, change the transaction state (other than by using a subtransaction for error recovery), access sequences (e.g., by calling `currval`) or make persistent changes to settings. They should be labeled parallel restricted if they access temporary tables, client connection state, cursors, prepared statements, or miscellaneous backend-local state which the system cannot synchronize in parallel mode (e.g., `setseed` cannot be executed other than by the group leader because a change made by another process would not be reflected in the leader). In general, if a function is labeled as being safe when it is restricted or unsafe, or if it is labeled as being restricted when it is in fact unsafe, it may throw errors or produce wrong answers when used in a parallel query. C-language functions could in theory exhibit totally undefined behavior if mislabeled, since there is no way for the system to protect itself against arbitrary C code, but in most likely cases the result will be no worse than for any other function. If in doubt, functions should be labeled as `UNSAFE`, which is the default.

`COST execution_cost`

A positive number giving the estimated execution cost for the function, in units of [cpu_operator_cost](#). If the function returns a set, this is the cost per returned row. If the cost is not specified, 1 unit is assumed for C-language and internal functions, and 100 units for functions in all other languages. Larger values cause the planner to try to avoid evaluating the function more often than necessary.

`ROWS result_rows`

A positive number giving the estimated number of rows that the planner should expect the function to return. This is only allowed when the function is declared to return a set. The default assumption is 1000 rows.

SUPPORT *support_function*

The name (optionally schema-qualified) of a *planner support function* to use for this function. See [Section 39.11](#) for details. You must be superuser to use this option.

configuration_parameter
value

The SET clause causes the specified configuration parameter to be set to the specified value when the function is entered, and then restored to its prior value when the function exits. SET FROM CURRENT saves the value of the parameter that is current when CREATE FUNCTION is executed as the value to be applied when the function is entered.

If a SET clause is attached to a function, then the effects of a SET LOCAL command executed inside the function for the same variable are restricted to the function: the configuration parameter's prior value is still restored at function exit. However, an ordinary SET command (without LOCAL) overrides the SET clause, much as it would do for a previous SET LOCAL command: the effects of such a command will persist after function exit, unless the current transaction is rolled back.

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

definition

A string constant defining the function; the meaning depends on the language. It can be an internal function name, the path to an object file, an SQL command, or text in a procedural language.

It is often helpful to use dollar quoting (see [Section 4.1.2.4](#)) to write the function definition string, rather than the normal single quote syntax. Without dollar quoting, any single quotes or backslashes in the function definition must be escaped by doubling them.

obj_file, link_symbol

This form of the AS clause is used for dynamically loadable C language functions when the function name in the C language source code is not the same as the name of the SQL function. The string *obj_file* is the name of the shared library file containing the compiled C function, and is interpreted as for the [LOAD](#) command. The string *link_symbol* is the function's link symbol, that is, the name of the function in the C language source code. If the link symbol is omitted, it is assumed to be the same as the name of the SQL function being defined. The C names of all functions must be different, so you must give overloaded C functions different C names (for example, use the argument types as part of the C names).

When repeated CREATE FUNCTION calls refer to the same object file, the file is only loaded once per session. To unload and reload the file (perhaps during development), start a new session.

sql_body

The body of a LANGUAGE SQL function. This can either be a single statement

RETURN *expression*

or a block

```
BEGIN ATOMIC
    statement;
    statement;
    ...
    statement;
END
```

This is similar to writing the text of the function body as a string constant (see *definition* above), but there are some differences: This form only works for LANGUAGE SQL, the string constant form works for all languages. This form is parsed at function definition time, the string constant form is parsed at execution time; therefore this form cannot support polymorphic argument types and other

constructs that are not resolvable at function definition time. This form tracks dependencies between the function and objects used in the function body, so `DROP ... CASCADE` will work correctly, whereas the form using string literals may leave dangling functions. Finally, this form is more compatible with the SQL standard and other SQL implementations.

Overloading

Postgres Pro allows function *overloading*; that is, the same name can be used for several different functions so long as they have distinct input argument types. Whether or not you use it, this capability entails security precautions when calling functions in databases where some users mistrust other users; see [Section 10.3](#).

Two functions are considered the same if they have the same names and *input* argument types, ignoring any `OUT` parameters. Thus for example these declarations conflict:

```
CREATE FUNCTION foo(int) ...
CREATE FUNCTION foo(int, out text) ...
```

Functions that have different argument type lists will not be considered to conflict at creation time, but if defaults are provided they might conflict in use. For example, consider

```
CREATE FUNCTION foo(int) ...
CREATE FUNCTION foo(int, int default 42) ...
```

A call `foo(10)` will fail due to the ambiguity about which function should be called.

Notes

The full SQL type syntax is allowed for declaring a function's arguments and return value. However, parenthesized type modifiers (e.g., the precision field for type `numeric`) are discarded by `CREATE FUNCTION`. Thus for example `CREATE FUNCTION foo (varchar(10)) ...` is exactly the same as `CREATE FUNCTION foo (varchar)`

When replacing an existing function with `CREATE OR REPLACE FUNCTION`, there are restrictions on changing parameter names. You cannot change the name already assigned to any input parameter (although you can add names to parameters that had none before). If there is more than one output parameter, you cannot change the names of the output parameters, because that would change the column names of the anonymous composite type that describes the function's result. These restrictions are made to ensure that existing calls of the function do not stop working when it is replaced.

If a function is declared `STRICT` with a `VARIADIC` argument, the strictness check tests that the variadic array *as a whole* is non-null. The function will still be called if the array has null elements.

Examples

Add two integers using an SQL function:

```
CREATE FUNCTION add(integer, integer) RETURNS integer
  AS 'select $1 + $2;'
  LANGUAGE SQL
  IMMUTABLE
  RETURNS NULL ON NULL INPUT;
```

The same function written in a more SQL-conforming style, using argument names and an unquoted body:

```
CREATE FUNCTION add(a integer, b integer) RETURNS integer
  LANGUAGE SQL
  IMMUTABLE
  RETURNS NULL ON NULL INPUT
  RETURN a + b;
```

Increment an integer, making use of an argument name, in PL/pgSQL:

```
CREATE OR REPLACE FUNCTION increment(i integer) RETURNS integer AS $$
    BEGIN
        RETURN i + 1;
    END;
$$ LANGUAGE plpgsql;
```

Return a record containing multiple output parameters:

```
CREATE FUNCTION dup(in int, out f1 int, out f2 text)
    AS $$ SELECT $1, CAST($1 AS text) || ' is text' $$
    LANGUAGE SQL;

SELECT * FROM dup(42);
```

You can do the same thing more verbosely with an explicitly named composite type:

```
CREATE TYPE dup_result AS (f1 int, f2 text);

CREATE FUNCTION dup(int) RETURNS dup_result
    AS $$ SELECT $1, CAST($1 AS text) || ' is text' $$
    LANGUAGE SQL;

SELECT * FROM dup(42);
```

Another way to return multiple columns is to use a TABLE function:

```
CREATE FUNCTION dup(int) RETURNS TABLE(f1 int, f2 text)
    AS $$ SELECT $1, CAST($1 AS text) || ' is text' $$
    LANGUAGE SQL;

SELECT * FROM dup(42);
```

However, a TABLE function is different from the preceding examples, because it actually returns a *set* of records, not just one record.

Writing SECURITY DEFINER Functions Safely

Because a SECURITY DEFINER function is executed with the privileges of the user that owns it, care is needed to ensure that the function cannot be misused. For security, [search_path](#) should be set to exclude any schemas writable by untrusted users. This prevents malicious users from creating objects (e.g., tables, functions, and operators) that mask objects intended to be used by the function. Particularly important in this regard is the temporary-table schema, which is searched first by default, and is normally writable by anyone. A secure arrangement can be obtained by forcing the temporary schema to be searched last. To do this, write `pg_temp` as the last entry in `search_path`. This function illustrates safe usage:

```
CREATE FUNCTION check_password(uname TEXT, pass TEXT)
    RETURNS BOOLEAN AS $$
    DECLARE passed BOOLEAN;
    BEGIN
        SELECT (pwd = $2) INTO passed
        FROM   pwds
        WHERE  username = $1;

        RETURN passed;
    END;
$$ LANGUAGE plpgsql
    SECURITY DEFINER
    -- Set a secure search_path: trusted schema(s), then 'pg_temp'.
    SET search_path = admin, pg_temp;
```

This function's intention is to access a table `admin.pwds`. But without the SET clause, or with a SET clause mentioning only `admin`, the function could be subverted by creating a temporary table named `pwds`.

If the security definer function intends to create roles, and if it is running as a non-superuser, `create_role_self_grant` should also be set to a known value using the `SET` clause.

Another point to keep in mind is that by default, execute privilege is granted to `PUBLIC` for newly created functions (see [Section 5.8](#) for more information). Frequently you will wish to restrict use of a security definer function to only some users. To do that, you must revoke the default `PUBLIC` privileges and then grant execute privilege selectively. To avoid having a window where the new function is accessible to all, create it and set the privileges within a single transaction. For example:

```
BEGIN;
CREATE FUNCTION check_password(uname TEXT, pass TEXT) ... SECURITY DEFINER;
REVOKE ALL ON FUNCTION check_password(uname TEXT, pass TEXT) FROM PUBLIC;
GRANT EXECUTE ON FUNCTION check_password(uname TEXT, pass TEXT) TO admins;
COMMIT;
```

Compatibility

A `CREATE FUNCTION` command is defined in the SQL standard. The Postgres Pro implementation can be used in a compatible way but has many extensions. Conversely, the SQL standard specifies a number of optional features that are not implemented in Postgres Pro.

The following are important compatibility issues:

- `OR REPLACE` is a Postgres Pro extension.
- For compatibility with some other database systems, *argmode* can be written either before or after *argname*. But only the first way is standard-compliant.
- For parameter defaults, the SQL standard specifies only the syntax with the `DEFAULT` key word. The syntax with `=` is used in T-SQL and Firebird.
- The `SETOF` modifier is a Postgres Pro extension.
- Only `SQL` is standardized as a language.
- All other attributes except `CALLED ON NULL INPUT` and `RETURNS NULL ON NULL INPUT` are not standardized.
- For the body of `LANGUAGE SQL` functions, the SQL standard only specifies the *sql_body* form.

Simple `LANGUAGE SQL` functions can be written in a way that is both standard-conforming and portable to other implementations. More complex functions using advanced features, optimization attributes, or other languages will necessarily be specific to Postgres Pro in a significant way.

See Also

[ALTER FUNCTION](#), [DROP FUNCTION](#), [GRANT](#), [LOAD](#), [REVOKE](#)

CREATE GROUP

CREATE GROUP — define a new database role

Synopsis

```
CREATE GROUP name [ [ WITH ] option [ ... ] ]
```

where *option* can be:

```
SUPERUSER | NOSUPERUSER
| CREATEDB | NOCREATEDB
| CREATEROLE | NOCREATEROLE
| INHERIT | NOINHERIT
| LOGIN | NOLOGIN
| REPLICATION | NOREPLICATION
| BYPASSRLS | NOBYPASSRLS
| CONNECTION LIMIT connlimit
| [ ENCRYPTED ] PASSWORD 'password' | PASSWORD NULL
| VALID UNTIL 'timestamp'
| IN ROLE role_name [, ...]
| IN GROUP role_name [, ...]
| ROLE role_name [, ...]
| ADMIN role_name [, ...]
| USER role_name [, ...]
| SYSID uid
```

Description

CREATE GROUP is now an alias for [CREATE ROLE](#).

Compatibility

There is no CREATE GROUP statement in the SQL standard.

See Also

[CREATE ROLE](#)

CREATE INDEX

CREATE INDEX — define a new index

Synopsis

```
CREATE [ UNIQUE ] INDEX [ CONCURRENTLY ] [ [ IF NOT EXISTS ] name ] ON
[ ONLY ] table_name [ USING method ]
    ( { column_name | ( expression ) } [ COLLATE collation ] [ opclass
[ ( opclass_parameter = value [, ... ] ) ] ] [ ASC | DESC ] [ NULLS { FIRST | LAST } ]
[, ...] )
    [ INCLUDE ( column_name [, ...] ) ]
    [ NULLS [ NOT ] DISTINCT ]
    [ WITH ( storage_parameter [= value] [, ... ] ) ]
    [ TABLESPACE tablespace_name ]
    [ WHERE predicate ]
```

Description

`CREATE INDEX` constructs an index on the specified column(s) of the specified relation, which can be a table or a materialized view. Indexes are primarily used to enhance database performance (though inappropriate use can result in slower performance).

The key field(s) for the index are specified as column names, or alternatively as expressions written in parentheses. Multiple fields can be specified if the index method supports multicolumn indexes.

An index field can be an expression computed from the values of one or more columns of the table row. This feature can be used to obtain fast access to data based on some transformation of the basic data. For example, an index computed on `upper(col)` would allow the clause `WHERE upper(col) = 'JIM'` to use an index.

Postgres Pro provides the index methods B-tree, hash, GiST, SP-GiST, GIN, RUM, and BRIN. Users can also define their own index methods, but that is fairly complicated.

When the `WHERE` clause is present, a *partial index* is created. A partial index is an index that contains entries for only a portion of a table, usually a portion that is more useful for indexing than the rest of the table. For example, if you have a table that contains both billed and unbilled orders where the unbilled orders take up a small fraction of the total table and yet that is an often used section, you can improve performance by creating an index on just that portion. Another possible application is to use `WHERE` with `UNIQUE` to enforce uniqueness over a subset of a table. See [Section 11.8](#) for more discussion.

The expression used in the `WHERE` clause can refer only to columns of the underlying table, but it can use all columns, not just the ones being indexed. Presently, subqueries and aggregate expressions are also forbidden in `WHERE`. The same restrictions apply to index fields that are expressions.

All functions and operators used in an index definition must be “immutable”, that is, their results must depend only on their arguments and never on any outside influence (such as the contents of another table or the current time). This restriction ensures that the behavior of the index is well-defined. To use a user-defined function in an index expression or `WHERE` clause, remember to mark the function immutable when you create it.

Parameters

`UNIQUE`

Causes the system to check for duplicate values in the table when the index is created (if data already exist) and each time data is added. Attempts to insert or update data which would result in duplicate entries will generate an error.

Additional restrictions apply when unique indexes are applied to partitioned tables; see [CREATE TABLE](#).

CONCURRENTLY

When this option is used, Postgres Pro will build the index without taking any locks that prevent concurrent inserts, updates, or deletes on the table; whereas a standard index build locks out writes (but not reads) on the table until it's done. There are several caveats to be aware of when using this option — see [Building Indexes Concurrently](#) below.

For temporary tables, `CREATE INDEX` is always non-concurrent, as no other session can access them, and non-concurrent index creation is cheaper.

IF NOT EXISTS

Do not throw an error if a relation with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing index is anything like the one that would have been created. Index name is required when `IF NOT EXISTS` is specified.

INCLUDE

The optional `INCLUDE` clause specifies a list of columns which will be included in the index as *non-key* columns. A non-key column cannot be used in an index scan search qualification, and it is disregarded for purposes of any uniqueness or exclusion constraint enforced by the index. However, an index-only scan can return the contents of non-key columns without having to visit the index's table, since they are available directly from the index entry. Thus, addition of non-key columns allows index-only scans to be used for queries that otherwise could not use them.

It's wise to be conservative about adding non-key columns to an index, especially wide columns. If an index tuple exceeds the maximum size allowed for the index type, data insertion will fail. In any case, non-key columns duplicate data from the index's table and bloat the size of the index, thus potentially slowing searches. Furthermore, B-tree deduplication is never used with indexes that have a non-key column.

Columns listed in the `INCLUDE` clause don't need appropriate operator classes; the clause can include columns whose data types don't have operator classes defined for a given access method.

Expressions are not supported as included columns since they cannot be used in index-only scans.

Currently, the B-tree, GiST and SP-GiST index access methods support this feature. In these indexes, the values of columns listed in the `INCLUDE` clause are included in leaf tuples which correspond to heap tuples, but are not included in upper-level index entries used for tree navigation.

name

The name of the index to be created. No schema name can be included here; the index is always created in the same schema as its parent table. The name of the index must be distinct from the name of any other relation (table, sequence, index, view, materialized view, or foreign table) in that schema. If the name is omitted, Postgres Pro chooses a suitable name based on the parent table's name and the indexed column name(s).

ONLY

Indicates not to recurse creating indexes on partitions, if the table is partitioned. The default is to recurse.

table_name

The name (possibly schema-qualified) of the table to be indexed.

method

The name of the index method to be used. Choices are `btree`, `hash`, `gist`, `spgist`, `gin`, `brin`, or user-installed access methods . The default method is `btree`.

column_name

The name of a column of the table.

expression

An expression based on one or more columns of the table. The expression usually must be written with surrounding parentheses, as shown in the syntax. However, the parentheses can be omitted if the expression has the form of a function call.

collation

The name of the collation to use for the index. By default, the index uses the collation declared for the column to be indexed or the result collation of the expression to be indexed. Indexes with non-default collations can be useful for queries that involve expressions using non-default collations.

opclass

The name of an operator class. See below for details.

opclass_parameter

The name of an operator class parameter. See below for details.

ASC

Specifies ascending sort order (which is the default).

DESC

Specifies descending sort order.

NULLS FIRST

Specifies that nulls sort before non-nulls. This is the default when DESC is specified.

NULLS LAST

Specifies that nulls sort after non-nulls. This is the default when DESC is not specified.

NULLS DISTINCT

NULLS NOT DISTINCT

Specifies whether for a unique index, null values should be considered distinct (not equal). The default is that they are distinct, so that a unique index could contain multiple null values in a column.

storage_parameter

The name of an index-method-specific storage parameter. See [Index Storage Parameters](#) below for details.

tablespace_name

The tablespace in which to create the index. If not specified, [default_tablespace](#) is consulted, or [temp_tablespaces](#) for indexes on temporary tables.

predicate

The constraint expression for a partial index.

Index Storage Parameters

The optional WITH clause specifies *storage parameters* for the index. Each index method has its own set of allowed storage parameters. The B-tree, hash, GiST and SP-GiST index methods all accept this parameter:

`fillfactor (integer)`

The fillfactor for an index is a percentage that determines how full the index method will try to pack index pages. For B-trees, leaf pages are filled to this percentage during initial index builds, and also when extending the index at the right (adding new largest key values). If pages subsequently become completely full, they will be split, leading to fragmentation of the on-disk index structure. B-trees use a default fillfactor of 90, but any integer value from 10 to 100 can be selected.

B-tree indexes on tables where many inserts and/or updates are anticipated can benefit from lower fillfactor settings at `CREATE INDEX` time (following bulk loading into the table). Values in the range of 50 - 90 can usefully “smooth out” the *rate* of page splits during the early life of the B-tree index (lowering fillfactor like this may even lower the absolute number of page splits, though this effect is highly workload dependent). The B-tree bottom-up index deletion technique described in [Section 65.1.4.2](#) is dependent on having some “extra” space on pages to store “extra” tuple versions, and so can be affected by fillfactor (though the effect is usually not significant).

In other specific cases it might be useful to increase fillfactor to 100 at `CREATE INDEX` time as a way of maximizing space utilization. You should only consider this when you are completely sure that the table is static (i.e. that it will never be affected by either inserts or updates). A fillfactor setting of 100 otherwise risks *harming* performance: even a few updates or inserts will cause a sudden flood of page splits.

The other index methods use fillfactor in different but roughly analogous ways; the default fillfactor varies between methods.

B-tree indexes additionally accept this parameter:

`deduplicate_items (boolean)`

Controls usage of the B-tree deduplication technique described in [Section 65.1.4.3](#). Set to `ON` or `OFF` to enable or disable the optimization. (Alternative spellings of `ON` and `OFF` are allowed as described in [Section 19.1](#).) The default is `ON`.

Note

Turning `deduplicate_items` off via `ALTER INDEX` prevents future insertions from triggering deduplication, but does not in itself make existing posting list tuples use the standard tuple representation.

GiST indexes additionally accept this parameter:

`buffering (enum)`

Determines whether the buffered build technique described in [Section 65.2.4.1](#) is used to build the index. With `OFF` buffering is disabled, with `ON` it is enabled, and with `AUTO` it is initially disabled, but is turned on on-the-fly once the index size reaches `effective_cache_size`. The default is `AUTO`. Note that if sorted build is possible, it will be used instead of buffered build unless `buffering=ON` is specified.

GIN indexes accept different parameters:

`fastupdate (boolean)`

This setting controls usage of the fast update technique described in [Section 65.4.4.1](#). It is a Boolean parameter: `ON` enables fast update, `OFF` disables it. The default is `ON`.

Note

Turning `fastupdate` off via `ALTER INDEX` prevents future insertions from going into the list of pending index entries, but does not in itself flush previous entries. You might want to `VACUUM`

the table or call `gin_clean_pending_list` function afterward to ensure the pending list is emptied.

`gin_pending_list_limit (integer)`

Custom [gin_pending_list_limit](#) parameter. This value is specified in kilobytes.

RUM indexes accept different parameters:

`attach (text)`

This setting specifies the name of the column to attach as additional information.

`to (text)`

This setting specifies the name of the column to add ordering by.

`order_by_attach (boolean)`

This setting controls whether to use the `(attach, itempointer)` order instead of just `itempointer`.

BRIN indexes accept different parameters:

`pages_per_range (integer)`

Defines the number of table blocks that make up one block range for each entry of a BRIN index (see [Section 65.5.1](#) for more details). The default is `128`.

`autosummarize (boolean)`

Defines whether a summarization run is queued for the previous page range whenever an insertion is detected on the next one. See [Section 65.5.1.1](#) for more details. The default is `off`.

Building Indexes Concurrently

Creating an index can interfere with regular operation of a database. Normally Postgres Pro locks the table to be indexed against writes and performs the entire index build with a single scan of the table. Other transactions can still read the table, but if they try to insert, update, or delete rows in the table they will block until the index build is finished. This could have a severe effect if the system is a live production database. Very large tables can take many hours to be indexed, and even for smaller tables, an index build can lock out writers for periods that are unacceptably long for a production system.

Postgres Pro supports building indexes without locking out writes. This method is invoked by specifying the `CONCURRENTLY` option of `CREATE INDEX`. When this option is used, Postgres Pro must perform two scans of the table, and in addition it must wait for all existing transactions that could potentially modify or use the index to terminate. Thus this method requires more total work than a standard index build and takes significantly longer to complete. However, since it allows normal operations to continue while the index is built, this method is useful for adding new indexes in a production environment. Of course, the extra CPU and I/O load imposed by the index creation might slow other operations.

In a concurrent index build, the index is actually entered as an “invalid” index into the system catalogs in one transaction, then two table scans occur in two more transactions. Before each table scan, the index build must wait for existing transactions that have modified the table to terminate. After the second scan, the index build must wait for any transactions that have a snapshot (see [Chapter 13](#)) predating the second scan to terminate, including transactions used by any phase of concurrent index builds on other tables, if the indexes involved are partial or have columns that are not simple column references. Then finally the index can be marked “valid” and ready for use, and the `CREATE INDEX` command terminates. Even then, however, the index may not be immediately usable for queries: in the worst case, it cannot be used as long as transactions exist that predate the start of the index build.

If a problem arises while scanning the table, such as a deadlock or a uniqueness violation in a unique index, the `CREATE INDEX` command will fail but leave behind an “invalid” index. If this happens while build an index concurrently on a partitioned table, the command can also leave behind “valid” or “invalid”

indexes on table partitions. The invalid index will be ignored for querying purposes because it might be incomplete; however it will still consume update overhead. The `psql \d` command will report such an index as `INVALID`:

```
postgres=# \d tab
          Table "public.tab"
  Column | Type   | Collation | Nullable | Default
-----+-----+-----+-----+-----
 col    | integer |           |          |
Indexes:
    "idx" btree (col) INVALID
```

The recommended recovery method in such cases is to drop the index and try again to perform `CREATE INDEX CONCURRENTLY`. (Another possibility is to rebuild the index with `REINDEX INDEX CONCURRENTLY`).

Another caveat when building a unique index concurrently is that the uniqueness constraint is already being enforced against other transactions when the second table scan begins. This means that constraint violations could be reported in other queries prior to the index becoming available for use, or even in cases where the index build eventually fails. Also, if a failure does occur in the second scan, the “invalid” index continues to enforce its uniqueness constraint afterwards.

Concurrent builds of expression indexes and partial indexes are supported. Errors occurring in the evaluation of these expressions could cause behavior similar to that described above for unique constraint violations.

Regular index builds permit other regular index builds on the same table to occur simultaneously, but only one concurrent index build can occur on a table at a time. In either case, schema modification of the table is not allowed while the index is being built. Another difference is that a regular `CREATE INDEX` command can be performed within a transaction block, but `CREATE INDEX CONCURRENTLY` cannot.

Notes

See [Chapter 11](#) for information about when indexes can be used, when they are not used, and in which particular situations they can be useful.

Currently, only the B-tree, GiST, GIN, and BRIN index methods support multiple-key-column indexes. Whether there can be multiple key columns is independent of whether `INCLUDE` columns can be added to the index. Indexes can have up to 32 columns, including `INCLUDE` columns. (This limit can be altered when building Postgres Pro.) Only B-tree currently supports unique indexes.

An *operator class* with optional parameters can be specified for each column of an index. The operator class identifies the operators to be used by the index for that column. For example, a B-tree index on four-byte integers would use the `int4_ops` class; this operator class includes comparison functions for four-byte integers. In practice the default operator class for the column's data type is usually sufficient. The main point of having operator classes is that for some data types, there could be more than one meaningful ordering. For example, we might want to sort a complex-number data type either by absolute value or by real part. We could do this by defining two operator classes for the data type and then selecting the proper class when creating an index. More information about operator classes is in [Section 11.10](#) and in [Section 39.16](#).

When `CREATE INDEX` is invoked on a partitioned table, the default behavior is to recurse to all partitions to ensure they all have matching indexes. Each partition is first checked to determine whether an equivalent index already exists, and if so, that index will become attached as a partition index to the index being created, which will become its parent index. If no matching index exists, a new index will be created and automatically attached; the name of the new index in each partition will be determined as if no index name had been specified in the command. If the `ONLY` option is specified, no recursion is done, and the index is marked invalid. (`ALTER INDEX ... ATTACH PARTITION` marks the index valid, once all partitions acquire matching indexes.) Note, however, that any partition that is created in the future using `CREATE TABLE ... PARTITION OF` will automatically have a matching index, regardless of whether `ONLY` is specified.

For index methods that support ordered scans (currently, only B-tree), the optional clauses `ASC`, `DESC`, `NULLS FIRST`, and/or `NULLS LAST` can be specified to modify the sort ordering of the index. Since an ordered index can be scanned either forward or backward, it is not normally useful to create a single-column `DESC` index — that sort ordering is already available with a regular index. The value of these options is that multicolumn indexes can be created that match the sort ordering requested by a mixed-ordering query, such as `SELECT ... ORDER BY x ASC, y DESC`. The `NULLS` options are useful if you need to support “nulls sort low” behavior, rather than the default “nulls sort high”, in queries that depend on indexes to avoid sorting steps.

The system regularly collects statistics on all of a table's columns. Newly-created non-expression indexes can immediately use these statistics to determine an index's usefulness. For new expression indexes, it is necessary to run `ANALYZE` or wait for the `autovacuum daemon` to analyze the table to generate statistics for these indexes.

While `CREATE INDEX` is running, the `search_path` is temporarily changed to `pg_catalog, pg_temp`.

For most index methods, the speed of creating an index is dependent on the setting of `maintenance_work_mem`. Larger values will reduce the time needed for index creation, so long as you don't make it larger than the amount of memory really available, which would drive the machine into swapping.

Postgres Pro can build indexes while leveraging multiple CPUs in order to process the table rows faster. This feature is known as *parallel index build*. For index methods that support building indexes in parallel (currently, B-tree and BRIN), `maintenance_work_mem` specifies the maximum amount of memory that can be used by each index build operation as a whole, regardless of how many worker processes were started. Generally, a cost model automatically determines how many worker processes should be requested, if any.

Parallel index builds may benefit from increasing `maintenance_work_mem` where an equivalent serial index build will see little or no benefit. Note that `maintenance_work_mem` may influence the number of worker processes requested, since parallel workers must have at least a 32MB share of the total `maintenance_work_mem` budget. There must also be a remaining 32MB share for the leader process. Increasing `max_parallel_maintenance_workers` may allow more workers to be used, which will reduce the time needed for index creation, so long as the index build is not already I/O bound. Of course, there should also be sufficient CPU capacity that would otherwise lie idle.

Setting a value for `parallel_workers` via `ALTER TABLE` directly controls how many parallel worker processes will be requested by a `CREATE INDEX` against the table. This bypasses the cost model completely, and prevents `maintenance_work_mem` from affecting how many parallel workers are requested. Setting `parallel_workers` to 0 via `ALTER TABLE` will disable parallel index builds on the table in all cases.

Tip

You might want to reset `parallel_workers` after setting it as part of tuning an index build. This avoids inadvertent changes to query plans, since `parallel_workers` affects *all* parallel table scans.

While `CREATE INDEX` with the `CONCURRENTLY` option supports parallel builds without special restrictions, only the first table scan is actually performed in parallel.

Use `DROP INDEX` to remove an index.

Like any long-running transaction, `CREATE INDEX` on a table can affect which tuples can be removed by concurrent `VACUUM` on any other table.

Prior releases of Postgres Pro also had an R-tree index method. This method has been removed because it had no significant advantages over the GiST method. If `USING rtree` is specified, `CREATE INDEX` will interpret it as `USING gist`, to simplify conversion of old databases to GiST.

Each backend running `CREATE INDEX` will report its progress in the `pg_stat_progress_create_index` view. See [Section 27.4.4](#) for details.

Examples

To create a unique B-tree index on the column `title` in the table `films`:

```
CREATE UNIQUE INDEX title_idx ON films (title);
```

To create a unique B-tree index on the column `title` with included columns `director` and `rating` in the table `films`:

```
CREATE UNIQUE INDEX title_idx ON films (title) INCLUDE (director, rating);
```

To create a B-Tree index with deduplication disabled:

```
CREATE INDEX title_idx ON films (title) WITH (deduplicate_items = off);
```

To create an index on the expression `lower(title)`, allowing efficient case-insensitive searches:

```
CREATE INDEX ON films ((lower(title)));
```

(In this example we have chosen to omit the index name, so the system will choose a name, typically `films_lower_idx`.)

To create an index with non-default collation:

```
CREATE INDEX title_idx_german ON films (title COLLATE "de_DE");
```

To create an index with non-default sort ordering of nulls:

```
CREATE INDEX title_idx_nulls_low ON films (title NULLS FIRST);
```

To create an index with non-default fill factor:

```
CREATE UNIQUE INDEX title_idx ON films (title) WITH (fillfactor = 70);
```

To create a GIN index with fast updates disabled:

```
CREATE INDEX gin_idx ON documents_table USING GIN (locations) WITH (fastupdate = off);
```

To create an index on the column `code` in the table `films` and have the index reside in the tablespace `indexspace`:

```
CREATE INDEX code_idx ON films (code) TABLESPACE indexspace;
```

To create a GiST index on a point attribute so that we can efficiently use box operators on the result of the conversion function:

```
CREATE INDEX pointloc
  ON points USING gist (box(location,location));
SELECT * FROM points
  WHERE box(location,location) && '(0,0),(1,1)::box;
```

To create an index without locking out writes to the table:

```
CREATE INDEX CONCURRENTLY sales_quantity_index ON sales_table (quantity);
```

Compatibility

`CREATE INDEX` is a Postgres Pro language extension. There are no provisions for indexes in the SQL standard.

See Also

[ALTER INDEX](#), [DROP INDEX](#), [REINDEX](#), [Section 27.4.4](#)

CREATE LANGUAGE

CREATE LANGUAGE — define a new procedural language

Synopsis

```
CREATE [ OR REPLACE ] [ TRUSTED ] [ PROCEDURAL ] LANGUAGE name
    HANDLER call_handler [ INLINE inline_handler ] [ VALIDATOR valfunction ]
CREATE [ OR REPLACE ] [ TRUSTED ] [ PROCEDURAL ] LANGUAGE name
```

Description

CREATE LANGUAGE registers a new procedural language with a Postgres Pro database. Subsequently, functions and procedures can be defined in this new language.

CREATE LANGUAGE effectively associates the language name with handler function(s) that are responsible for executing functions written in the language. Refer to [Chapter 57](#) for more information about language handlers.

CREATE OR REPLACE LANGUAGE will either create a new language, or replace an existing definition. If the language already exists, its parameters are updated according to the command, but the language's ownership and permissions settings do not change, and any existing functions written in the language are assumed to still be valid.

One must have the Postgres Pro superuser privilege to register a new language or change an existing language's parameters. However, once the language is created it is valid to assign ownership of it to a non-superuser, who may then drop it, change its permissions, rename it, or assign it to a new owner. (Do not, however, assign ownership of the underlying C functions to a non-superuser; that would create a privilege escalation path for that user.)

The form of CREATE LANGUAGE that does not supply any handler function is obsolete. For backwards compatibility with old dump files, it is interpreted as CREATE EXTENSION. That will work if the language has been packaged into an extension of the same name, which is the conventional way to set up procedural languages.

Parameters

TRUSTED

TRUSTED specifies that the language does not grant access to data that the user would not otherwise have. If this key word is omitted when registering the language, only users with the Postgres Pro superuser privilege can use this language to create new functions.

PROCEDURAL

This is a noise word.

name

The name of the new procedural language. The name must be unique among the languages in the database.

HANDLER *call_handler*

call_handler is the name of a previously registered function that will be called to execute the procedural language's functions. The call handler for a procedural language must be written in a compiled language such as C with version 1 call convention and registered with Postgres Pro as a function taking no arguments and returning the `language_handler` type, a placeholder type that is simply used to identify the function as a call handler.

`INLINE inline_handler`

inline_handler is the name of a previously registered function that will be called to execute an anonymous code block (`DO` command) in this language. If no *inline_handler* function is specified, the language does not support anonymous code blocks. The handler function must take one argument of type `internal`, which will be the `DO` command's internal representation, and it will typically return `void`. The return value of the handler is ignored.

`VALIDATOR valfunction`

valfunction is the name of a previously registered function that will be called when a new function in the language is created, to validate the new function. If no validator function is specified, then a new function will not be checked when it is created. The validator function must take one argument of type `oid`, which will be the OID of the to-be-created function, and will typically return `void`.

A validator function would typically inspect the function body for syntactical correctness, but it can also look at other properties of the function, for example if the language cannot handle certain argument types. To signal an error, the validator function should use the `ereport()` function. The return value of the function is ignored.

Notes

Use `DROP LANGUAGE` to drop procedural languages.

The system catalog `pg_language` (see [Section 54.29](#)) records information about the currently installed languages. Also, the `psql` command `\dL` lists the installed languages.

To create functions in a procedural language, a user must have the `USAGE` privilege for the language. By default, `USAGE` is granted to `PUBLIC` (i.e., everyone) for trusted languages. This can be revoked if desired.

Procedural languages are local to individual databases. However, a language can be installed into the `template1` database, which will cause it to be available automatically in all subsequently-created databases.

Examples

A minimal sequence for creating a new procedural language is:

```
CREATE FUNCTION plsample_call_handler() RETURNS language_handler
AS '$libdir/plsample'
LANGUAGE C;
CREATE LANGUAGE plsample
HANDLER plsample_call_handler;
```

Typically that would be written in an extension's creation script, and users would do this to install the extension:

```
CREATE EXTENSION plsample;
```

Compatibility

`CREATE LANGUAGE` is a Postgres Pro extension.

See Also

[ALTER LANGUAGE](#), [CREATE FUNCTION](#), [DROP LANGUAGE](#), [GRANT](#), [REVOKE](#)

CREATE MATERIALIZED VIEW

CREATE MATERIALIZED VIEW — define a new materialized view

Synopsis

```
CREATE MATERIALIZED VIEW [ IF NOT EXISTS ] table_name
    [ ( column_name [, ...] ) ]
    [ USING method ]
    [ WITH ( storage_parameter [= value] [, ... ] ) ]
    [ TABLESPACE tablespace_name ]
    AS query
    [ WITH [ NO ] DATA ]
```

Description

CREATE MATERIALIZED VIEW defines a materialized view of a query. The query is executed and used to populate the view at the time the command is issued (unless WITH NO DATA is used) and may be refreshed later using REFRESH MATERIALIZED VIEW.

CREATE MATERIALIZED VIEW is similar to CREATE TABLE AS, except that it also remembers the query used to initialize the view, so that it can be refreshed later upon demand. A materialized view has many of the same properties as a table, but there is no support for temporary materialized views.

CREATE MATERIALIZED VIEW requires CREATE privilege on the schema used for the materialized view.

Parameters

IF NOT EXISTS

Do not throw an error if a materialized view with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing materialized view is anything like the one that would have been created.

table_name

The name (optionally schema-qualified) of the materialized view to be created. The name must be distinct from the name of any other relation (table, sequence, index, view, materialized view, or foreign table) in the same schema.

column_name

The name of a column in the new materialized view. If column names are not provided, they are taken from the output column names of the query.

USING *method*

This optional clause specifies the table access method to use to store the contents for the new materialized view; the method needs be an access method of type TABLE. See [Chapter 62](#) for more information. If this option is not specified, the default table access method is chosen for the new materialized view. See [default_table_access_method](#) for more information.

WITH (*storage_parameter* [= *value*] [, ...])

This clause specifies optional storage parameters for the new materialized view; see [Storage Parameters](#) in the [CREATE TABLE](#) documentation for more information. All parameters supported for CREATE TABLE are also supported for CREATE MATERIALIZED VIEW. See [CREATE TABLE](#) for more information.

TABLESPACE *tablespace_name*

The *tablespace_name* is the name of the tablespace in which the new materialized view is to be created. If not specified, [default_tablespace](#) is consulted.

query

A [SELECT](#), [TABLE](#), or [VALUES](#) command. This query will run within a security-restricted operation; in particular, calls to functions that themselves create temporary tables will fail. Also, while the query is running, the [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`.

WITH [NO] DATA

This clause specifies whether or not the materialized view should be populated at creation time. If not, the materialized view will be flagged as unscannable and cannot be queried until `REFRESH MATERIALIZED VIEW` is used.

Compatibility

`CREATE MATERIALIZED VIEW` is a Postgres Pro extension.

See Also

[ALTER MATERIALIZED VIEW](#), [CREATE TABLE AS](#), [CREATE VIEW](#), [DROP MATERIALIZED VIEW](#), [REFRESH MATERIALIZED VIEW](#)

CREATE OPERATOR

CREATE OPERATOR — define a new operator

Synopsis

```
CREATE OPERATOR name (  
    {FUNCTION|PROCEDURE} = function_name  
    [, LEFTARG = left_type ] [, RIGHTARG = right_type ]  
    [, COMMUTATOR = com_op ] [, NEGATOR = neg_op ]  
    [, RESTRICT = res_proc ] [, JOIN = join_proc ]  
    [, HASHES ] [, MERGES ]  
)
```

Description

CREATE OPERATOR defines a new operator, *name*. The user who defines an operator becomes its owner. If a schema name is given then the operator is created in the specified schema. Otherwise it is created in the current schema.

The operator name is a sequence of up to NAMEDATALEN-1 (63 by default) characters from the following list:

+ - * / < > = ~ ! @ # % ^ & | ` ?

There are a few restrictions on your choice of name:

- -- and /* cannot appear anywhere in an operator name, since they will be taken as the start of a comment.
- A multicharacter operator name cannot end in + or -, unless the name also contains at least one of these characters:

~ ! @ # % ^ & | ` ?

For example, @- is an allowed operator name, but *- is not. This restriction allows Postgres Pro to parse SQL-compliant commands without requiring spaces between tokens.

- The symbol => is reserved by the SQL grammar, so it cannot be used as an operator name.

The operator != is mapped to <> on input, so these two names are always equivalent.

For binary operators, both LEFTARG and RIGHTARG must be defined. For prefix operators only RIGHTARG should be defined. The *function_name* function must have been previously defined using CREATE FUNCTION and must be defined to accept the correct number of arguments (either one or two) of the indicated types.

In the syntax of CREATE OPERATOR, the keywords FUNCTION and PROCEDURE are equivalent, but the referenced function must in any case be a function, not a procedure. The use of the keyword PROCEDURE here is historical and deprecated.

The other clauses specify optional operator optimization attributes. Their meaning is detailed in [Section 39.15](#).

To be able to create an operator, you must have USAGE privilege on the argument types and the return type, as well as EXECUTE privilege on the underlying function. If a commutator or negator operator is specified, you must own those operators.

Parameters

name

The name of the operator to be defined. See above for allowable characters. The name can be schema-qualified, for example CREATE OPERATOR myschema.+ (...). If not, then the operator is created in

the current schema. Two operators in the same schema can have the same name if they operate on different data types. This is called *overloading*.

function_name

The function used to implement this operator.

left_type

The data type of the operator's left operand, if any. This option would be omitted for a prefix operator.

right_type

The data type of the operator's right operand.

com_op

The commutator of this operator.

neg_op

The negator of this operator.

res_proc

The restriction selectivity estimator function for this operator.

join_proc

The join selectivity estimator function for this operator.

HASHES

Indicates this operator can support a hash join.

MERGES

Indicates this operator can support a merge join.

To give a schema-qualified operator name in *com_op* or the other optional arguments, use the `OPERATOR()` syntax, for example:

```
COMMUTATOR = OPERATOR(myschema.===) ,
```

Notes

Refer to [Section 39.14](#) and [Section 39.15](#) for further information.

When you are defining a self-commutative operator, you just do it. When you are defining a pair of commutative operators, things are a little trickier: how can the first one to be defined refer to the other one, which you haven't defined yet? There are three solutions to this problem:

- One way is to omit the `COMMUTATOR` clause in the first operator that you define, and then provide one in the second operator's definition. Since Postgres Pro knows that commutative operators come in pairs, when it sees the second definition it will automatically go back and fill in the missing `COMMUTATOR` clause in the first definition.
- Another, more straightforward way is just to include `COMMUTATOR` clauses in both definitions. When Postgres Pro processes the first definition and realizes that `COMMUTATOR` refers to a nonexistent operator, the system will make a dummy entry for that operator in the system catalog. This dummy entry will have valid data only for the operator name, left and right operand types, and owner, since that's all that Postgres Pro can deduce at this point. The first operator's catalog entry will link to this dummy entry. Later, when you define the second operator, the system updates the dummy entry with the additional information from the second definition. If you try to use the dummy operator before it's been filled in, you'll just get an error message.

- Alternatively, both operators can be defined without `COMMUTATOR` clauses and then `ALTER OPERATOR` can be used to set their commutator links. It's sufficient to `ALTER` either one of the pair.

In all three cases, you must own both operators in order to mark them as commutators.

Pairs of negator operators can be defined using the same methods as for commutator pairs.

It is not possible to specify an operator's lexical precedence in `CREATE OPERATOR`, because the parser's precedence behavior is hard-wired. See [Section 4.1.6](#) for precedence details.

The obsolete options `SORT1`, `SORT2`, `LTCMP`, and `GTCMP` were formerly used to specify the names of sort operators associated with a merge-joinable operator. This is no longer necessary, since information about associated operators is found by looking at B-tree operator families instead. If one of these options is given, it is ignored except for implicitly setting `MERGES` true.

Use `DROP OPERATOR` to delete user-defined operators from a database. Use `ALTER OPERATOR` to modify operators in a database.

Examples

The following command defines a new operator, area-equality, for the data type `box`:

```
CREATE OPERATOR === (  
    LEFTARG = box,  
    RIGHTARG = box,  
    FUNCTION = area_equal_function,  
    COMMUTATOR = ==,  
    NEGATOR = !=,  
    RESTRICT = area_restriction_function,  
    JOIN = area_join_function,  
    HASHES, MERGES  
);
```

Compatibility

`CREATE OPERATOR` is a Postgres Pro extension. There are no provisions for user-defined operators in the SQL standard.

See Also

[ALTER OPERATOR](#), [CREATE OPERATOR CLASS](#), [DROP OPERATOR](#)

CREATE OPERATOR CLASS

CREATE OPERATOR CLASS — define a new operator class

Synopsis

```
CREATE OPERATOR CLASS name [ DEFAULT ] FOR TYPE data_type
  USING index_method [ FAMILY family_name ] AS
  { OPERATOR strategy_number operator_name [ ( op_type, op_type ) ] [ FOR SEARCH | FOR
ORDER BY sort_family_name ]
    | FUNCTION support_number [ ( op_type [ , op_type ] ) ] function_name
    ( argument_type [, ...] )
    | STORAGE storage_type
  } [, ... ]
```

Description

CREATE OPERATOR CLASS creates a new operator class. An operator class defines how a particular data type can be used with an index. The operator class specifies that certain operators will fill particular roles or “strategies” for this data type and this index method. The operator class also specifies the support functions to be used by the index method when the operator class is selected for an index column. All the operators and functions used by an operator class must be defined before the operator class can be created.

If a schema name is given then the operator class is created in the specified schema. Otherwise it is created in the current schema. Two operator classes in the same schema can have the same name only if they are for different index methods.

The user who defines an operator class becomes its owner. Presently, the creating user must be a superuser. (This restriction is made because an erroneous operator class definition could confuse or even crash the server.)

CREATE OPERATOR CLASS does not presently check whether the operator class definition includes all the operators and functions required by the index method, nor whether the operators and functions form a self-consistent set. It is the user's responsibility to define a valid operator class.

Related operator classes can be grouped into *operator families*. To add a new operator class to an existing family, specify the FAMILY option in CREATE OPERATOR CLASS. Without this option, the new class is placed into a family named the same as the new class (creating that family if it doesn't already exist).

Refer to [Section 39.16](#) for further information.

Parameters

name

The name of the operator class to be created. The name can be schema-qualified.

DEFAULT

If present, the operator class will become the default operator class for its data type. At most one operator class can be the default for a specific data type and index method.

data_type

The column data type that this operator class is for.

index_method

The name of the index method this operator class is for.

family_name

The name of the existing operator family to add this operator class to. If not specified, a family named the same as the operator class is used (creating it, if it doesn't already exist).

strategy_number

The index method's strategy number for an operator associated with the operator class.

operator_name

The name (optionally schema-qualified) of an operator associated with the operator class.

op_type

In an `OPERATOR` clause, the operand data type(s) of the operator, or `NONE` to signify a prefix operator. The operand data types can be omitted in the normal case where they are the same as the operator class's data type.

In a `FUNCTION` clause, the operand data type(s) the function is intended to support, if different from the input data type(s) of the function (for B-tree comparison functions and hash functions) or the class's data type (for B-tree sort support functions, B-tree equal image functions, and all functions in GiST, SP-GiST, GIN and BRIN operator classes). These defaults are correct, and so *op_type* need not be specified in `FUNCTION` clauses, except for the case of a B-tree sort support function that is meant to support cross-data-type comparisons.

sort_family_name

The name (optionally schema-qualified) of an existing `btree` operator family that describes the sort ordering associated with an ordering operator.

If neither `FOR SEARCH` nor `FOR ORDER BY` is specified, `FOR SEARCH` is the default.

support_number

The index method's support function number for a function associated with the operator class.

function_name

The name (optionally schema-qualified) of a function that is an index method support function for the operator class.

argument_type

The parameter data type(s) of the function.

storage_type

The data type actually stored in the index. Normally this is the same as the column data type, but some index methods (currently GiST, GIN, SP-GiST and BRIN) allow it to be different. The `STORAGE` clause must be omitted unless the index method allows a different type to be used. If the column *data_type* is specified as `anyarray`, the *storage_type* can be declared as `anyelement` to indicate that the index entries are members of the element type belonging to the actual array type that each particular index is created for.

The `OPERATOR`, `FUNCTION`, and `STORAGE` clauses can appear in any order.

Notes

Because the index machinery does not check access permissions on functions before using them, including a function or operator in an operator class is tantamount to granting public execute permission on it. This is usually not an issue for the sorts of functions that are useful in an operator class.

The operators should not be defined by SQL functions. An SQL function is likely to be inlined into the calling query, which will prevent the optimizer from recognizing that the query matches an index.

Before PostgreSQL 8.4, the `OPERATOR` clause could include a `RECHECK` option. This is no longer supported because whether an index operator is “lossy” is now determined on-the-fly at run time. This allows efficient handling of cases where an operator might or might not be lossy.

Examples

The following example command defines a GiST index operator class for the data type `_int4` (array of `int4`).

```
CREATE OPERATOR CLASS gist__int_ops
  DEFAULT FOR TYPE _int4 USING gist AS
    OPERATOR          3      &&,
    OPERATOR          6      = (anyarray, anyarray),
    OPERATOR          7      @>,
    OPERATOR          8      <@,
    OPERATOR          20     @@ (_int4, query_int),
    FUNCTION           1      g_int_consistent (internal, _int4, smallint, oid,
internal),
    FUNCTION           2      g_int_union (internal, internal),
    FUNCTION           3      g_int_compress (internal),
    FUNCTION           4      g_int_decompress (internal),
    FUNCTION           5      g_int_penalty (internal, internal, internal),
    FUNCTION           6      g_int_picksplit (internal, internal),
    FUNCTION           7      g_int_same (_int4, _int4, internal);
```

Compatibility

`CREATE OPERATOR CLASS` is a Postgres Pro extension. There is no `CREATE OPERATOR CLASS` statement in the SQL standard.

See Also

[ALTER OPERATOR CLASS](#), [DROP OPERATOR CLASS](#), [CREATE OPERATOR FAMILY](#), [ALTER OPERATOR FAMILY](#)

CREATE OPERATOR FAMILY

CREATE OPERATOR FAMILY — define a new operator family

Synopsis

```
CREATE OPERATOR FAMILY name USING index_method
```

Description

CREATE OPERATOR FAMILY creates a new operator family. An operator family defines a collection of related operator classes, and perhaps some additional operators and support functions that are compatible with these operator classes but not essential for the functioning of any individual index. (Operators and functions that are essential to indexes should be grouped within the relevant operator class, rather than being “loose” in the operator family. Typically, single-data-type operators are bound to operator classes, while cross-data-type operators can be loose in an operator family containing operator classes for both data types.)

The new operator family is initially empty. It should be populated by issuing subsequent CREATE OPERATOR CLASS commands to add contained operator classes, and optionally ALTER OPERATOR FAMILY commands to add “loose” operators and their corresponding support functions.

If a schema name is given then the operator family is created in the specified schema. Otherwise it is created in the current schema. Two operator families in the same schema can have the same name only if they are for different index methods.

The user who defines an operator family becomes its owner. Presently, the creating user must be a superuser. (This restriction is made because an erroneous operator family definition could confuse or even crash the server.)

Refer to [Section 39.16](#) for further information.

Parameters

name

The name of the operator family to be created. The name can be schema-qualified.

index_method

The name of the index method this operator family is for.

Compatibility

CREATE OPERATOR FAMILY is a Postgres Pro extension. There is no CREATE OPERATOR FAMILY statement in the SQL standard.

See Also

[ALTER OPERATOR FAMILY](#), [DROP OPERATOR FAMILY](#), [CREATE OPERATOR CLASS](#), [ALTER OPERATOR CLASS](#), [DROP OPERATOR CLASS](#)

CREATE PACKAGE

CREATE PACKAGE — define a new package

Synopsis

```
CREATE [ OR REPLACE ] PACKAGE package_name package_element [ ... ]
```

Description

CREATE PACKAGE enters a new package into the current database.

A package is essentially a schema that helps to organize the named objects with a related purpose so it can be also created using [CREATE SCHEMA](#) and is subject to the same actions as a schema. A package should contain only functions, procedures and composite types.

CREATE OR REPLACE PACKAGE will either create a new package, or replace an existing definition. You can replace the existing package if it contains only functions and types, other objects must be previously dropped. If you replace the package without changing the function signature, its body is replaced, while dependent objects still remain. If the function signature changes, you would actually be creating a new, distinct function. The latter would only work when there are no dependent objects, otherwise the CREATE OR REPLACE PACKAGE is forced to fail. The same restriction applies to replacing types. When replacing packages, make sure that no other sessions use them (which can be done by stopping the application before replacing the packages).

Parameters

package_name

The name of a package to be created. The package name must be distinct from the name of any existing package or schema in the current database. The name cannot begin with `pg_`, as such names are reserved for system schemas.

package_element

An SQL statement defining an object to be created within the package. Currently, only CREATE FUNCTION, CREATE TYPE, and CREATE PROCEDURE are accepted as clauses within CREATE PACKAGE. The subcommands are treated essentially the same as separate commands issued after creating the package. By default, all the variables declared in the package initialization function, package functions and package procedures are public, so they can be called using dot notation from outside the package by other functions, procedures and anonymous blocks that import the package. The `#private` modifier defines functions and procedures as private, and the `#export` modifier defines which package variables are public. For more information on package functions and modifiers, see [Section 44.11](#).

Notes

To create a package, the invoking user must be a superuser or have the CREATE privilege for the current database.

Examples

Create a package named `counter` and create functions within it:

```
CREATE PACKAGE counter
CREATE FUNCTION __init__() RETURNS void AS $$ -- package initialization
#export on
DECLARE
    n int := 1; -- public package variable n
    k int := 3; -- public package variable k
```



```
BEGIN
    FOR i IN 1..10 LOOP
        n := n + n;
    END LOOP;
END;
$$

CREATE FUNCTION inc() RETURNS int AS $$ -- public package function inc()
BEGIN
    n := n + 1;
    RETURN n;
END;
$$
;
```

Notice that the individual subcommands do not end with semicolons, just as in [CREATE SCHEMA](#), and no language is specified when creating functions.

The `#export` on modifier inside the package initialization function renders all the declared variables public, meaning they are available outside the package. The same result is achieved if the package does not contain the `#export` modifier (default). The `inc()` function is public by default and is available outside the package, as any other package function or procedure.

The example below shows how the above-described package can be used.

```
DO $$
#import counter
BEGIN
    RAISE NOTICE '%', counter.n;
    RAISE NOTICE '%', counter.inc();
END;
$$;
```

```
NOTICE: 1024
NOTICE: 1025
```

A variation of the previous example:

```
CREATE PACKAGE foo
CREATE TYPE footype AS (a int, b int)

CREATE FUNCTION __init__() RETURNS void AS $$
#export y
DECLARE
    x int := 1; -- private package variable x
    y int := 5; -- public package variable y
BEGIN
    RAISE NOTICE 'foo initialized';
END;
$$

CREATE FUNCTION get_x() RETURNS int AS $$ -- public package function get_x()
BEGIN
    RETURN x;
END;
$$

CREATE PROCEDURE set_x(val int) AS $$ -- public package procedure set_x()
BEGIN
    x := val;
```

```

    CALL foo.check_x();
END;
$$

CREATE PROCEDURE check_x() AS $$ -- private package procedure check_x()
#private
BEGIN
    IF x <= 0 THEN
        RAISE INFO 'now x is not natural number';
    ELSE
        RAISE INFO 'now x is natural number';
    END IF;
END;
$$
;

```

In this example, the `y` variable, the `get_x()` function and the `set_x(int)` procedure of the `foo` package are public.

A variation of the previous example:

```

DO $$
#import foo
BEGIN
    RAISE INFO 'y = %', foo.y;
    RAISE INFO 'x = %', foo.get_x();
    CALL foo.set_x(25);
    RAISE INFO 'x = %', foo.get_x();
END;
$$;

INFO:  y = 5
INFO:  x = 1
INFO:  x set to natural value
INFO:  x = 25

```

The following is an equivalent way of achieving the same result using the `CREATE SCHEMA` command:

```

CREATE SCHEMA foo;
CREATE TYPE foo.footype AS (a int, b int);

CREATE FUNCTION foo.__init__() RETURNS void AS $$
#export y
DECLARE
    x int := 1; -- private package variable x
    y int := 5; -- public package variable y
BEGIN
    RAISE NOTICE 'foo initialized';
END;
$$ LANGUAGE plpgsql;

CREATE FUNCTION foo.get_x() RETURNS int AS $$ -- public package function get_x()
#package
BEGIN
    RETURN x;
END;
$$ LANGUAGE plpgsql;

CREATE PROCEDURE foo.set_x(val int) AS $$ -- public package procedure set_x()
#package

```

```
BEGIN
  x := val;
  CALL foo.check_x();
END;
$$ LANGUAGE plpgsql;

CREATE PROCEDURE foo.check_x() AS $$ -- private package procedure check_x()
#package
#private
BEGIN
  IF x <= 0 THEN
    RAISE INFO 'x set to not natural value';
  ELSE
    RAISE INFO 'x set to natural value';
  END IF;
END;
$$ LANGUAGE plpgsql;
```

Note that in this case you must use [#package modifiers](#).

See Also

[DROP PACKAGE](#)

CREATE POLICY

CREATE POLICY — define a new row-level security policy for a table

Synopsis

```
CREATE POLICY name ON table_name
[ AS { PERMISSIVE | RESTRICTIVE } ]
[ FOR { ALL | SELECT | INSERT | UPDATE | DELETE } ]
[ TO { role_name | PUBLIC | CURRENT_ROLE | CURRENT_USER | SESSION_USER } [, ...] ]
[ USING ( using_expression ) ]
[ WITH CHECK ( check_expression ) ]
```

Description

The `CREATE POLICY` command defines a new row-level security policy for a table. Note that row-level security must be enabled on the table (using `ALTER TABLE ... ENABLE ROW LEVEL SECURITY`) in order for created policies to be applied.

A policy grants the permission to select, insert, update, or delete rows that match the relevant policy expression. Existing table rows are checked against the expression specified in `USING`, while new rows that would be created via `INSERT` or `UPDATE` are checked against the expression specified in `WITH CHECK`. When a `USING` expression returns true for a given row then that row is visible to the user, while if false or null is returned then the row is not visible. When a `WITH CHECK` expression returns true for a row then that row is inserted or updated, while if false or null is returned then an error occurs.

For `INSERT`, `UPDATE`, and `MERGE` statements, `WITH CHECK` expressions are enforced after `BEFORE` triggers are fired, and before any actual data modifications are made. Thus a `BEFORE ROW` trigger may modify the data to be inserted, affecting the result of the security policy check. `WITH CHECK` expressions are enforced before any other constraints.

Policy names are per-table. Therefore, one policy name can be used for many different tables and have a definition for each table which is appropriate to that table.

Policies can be applied for specific commands or for specific roles. The default for newly created policies is that they apply for all commands and roles, unless otherwise specified. Multiple policies may apply to a single command; see below for more details. [Table 311](#) summarizes how the different types of policy apply to specific commands.

For policies that can have both `USING` and `WITH CHECK` expressions (`ALL` and `UPDATE`), if no `WITH CHECK` expression is defined, then the `USING` expression will be used both to determine which rows are visible (normal `USING` case) and which new rows will be allowed to be added (`WITH CHECK` case).

If row-level security is enabled for a table, but no applicable policies exist, a “default deny” policy is assumed, so that no rows will be visible or updatable.

Parameters

name

The name of the policy to be created. This must be distinct from the name of any other policy for the table.

table_name

The name (optionally schema-qualified) of the table the policy applies to.

`PERMISSIVE`

Specify that the policy is to be created as a permissive policy. All permissive policies which are applicable to a given query will be combined together using the Boolean “OR” operator. By creating

permissive policies, administrators can add to the set of records which can be accessed. Policies are permissive by default.

RESTRICTIVE

Specify that the policy is to be created as a restrictive policy. All restrictive policies which are applicable to a given query will be combined together using the Boolean “AND” operator. By creating restrictive policies, administrators can reduce the set of records which can be accessed as all restrictive policies must be passed for each record.

Note that there needs to be at least one permissive policy to grant access to records before restrictive policies can be usefully used to reduce that access. If only restrictive policies exist, then no records will be accessible. When a mix of permissive and restrictive policies are present, a record is only accessible if at least one of the permissive policies passes, in addition to all the restrictive policies.

command

The command to which the policy applies. Valid options are `ALL`, `SELECT`, `INSERT`, `UPDATE`, and `DELETE`. `ALL` is the default. See below for specifics regarding how these are applied.

role_name

The role(s) to which the policy is to be applied. The default is `PUBLIC`, which will apply the policy to all roles.

using_expression

Any SQL conditional expression (returning `boolean`). The conditional expression cannot contain any aggregate or window functions. This expression will be added to queries that refer to the table if row-level security is enabled. Rows for which the expression returns true will be visible. Any rows for which the expression returns false or null will not be visible to the user (in a `SELECT`), and will not be available for modification (in an `UPDATE` or `DELETE`). Such rows are silently suppressed; no error is reported.

check_expression

Any SQL conditional expression (returning `boolean`). The conditional expression cannot contain any aggregate or window functions. This expression will be used in `INSERT` and `UPDATE` queries against the table if row-level security is enabled. Only rows for which the expression evaluates to true will be allowed. An error will be thrown if the expression evaluates to false or null for any of the records inserted or any of the records that result from the update. Note that the *check_expression* is evaluated against the proposed new contents of the row, not the original contents.

Per-Command Policies

`ALL`

Using `ALL` for a policy means that it will apply to all commands, regardless of the type of command. If an `ALL` policy exists and more specific policies exist, then both the `ALL` policy and the more specific policy (or policies) will be applied. Additionally, `ALL` policies will be applied to both the selection side of a query and the modification side, using the `USING` expression for both cases if only a `USING` expression has been defined.

As an example, if an `UPDATE` is issued, then the `ALL` policy will be applicable both to what the `UPDATE` will be able to select as rows to be updated (applying the `USING` expression), and to the resulting updated rows, to check if they are permitted to be added to the table (applying the `WITH CHECK` expression, if defined, and the `USING` expression otherwise). If an `INSERT` or `UPDATE` command attempts to add rows to the table that do not pass the `ALL` policy's `WITH CHECK` expression, the entire command will be aborted.

`SELECT`

Using `SELECT` for a policy means that it will apply to `SELECT` queries and whenever `SELECT` permissions are required on the relation the policy is defined for. The result is that only those records from the

relation that pass the `SELECT` policy will be returned during a `SELECT` query, and that queries that require `SELECT` permissions, such as `UPDATE`, will also only see those records that are allowed by the `SELECT` policy. A `SELECT` policy cannot have a `WITH CHECK` expression, as it only applies in cases where records are being retrieved from the relation.

INSERT

Using `INSERT` for a policy means that it will apply to `INSERT` commands and `MERGE` commands that contain `INSERT` actions. Rows being inserted that do not pass this policy will result in a policy violation error, and the entire `INSERT` command will be aborted. An `INSERT` policy cannot have a `USING` expression, as it only applies in cases where records are being added to the relation.

Note that `INSERT` with `ON CONFLICT DO UPDATE` checks `INSERT` policies' `WITH CHECK` expressions only for rows appended to the relation by the `INSERT` path.

UPDATE

Using `UPDATE` for a policy means that it will apply to `UPDATE`, `SELECT FOR UPDATE` and `SELECT FOR SHARE` commands, as well as auxiliary `ON CONFLICT DO UPDATE` clauses of `INSERT` commands. `MERGE` commands containing `UPDATE` actions are affected as well. Since `UPDATE` involves pulling an existing record and replacing it with a new modified record, `UPDATE` policies accept both a `USING` expression and a `WITH CHECK` expression. The `USING` expression determines which records the `UPDATE` command will see to operate against, while the `WITH CHECK` expression defines which modified rows are allowed to be stored back into the relation.

Any rows whose updated values do not pass the `WITH CHECK` expression will cause an error, and the entire command will be aborted. If only a `USING` clause is specified, then that clause will be used for both `USING` and `WITH CHECK` cases.

Typically an `UPDATE` command also needs to read data from columns in the relation being updated (e.g., in a `WHERE` clause or a `RETURNING` clause, or in an expression on the right hand side of the `SET` clause). In this case, `SELECT` rights are also required on the relation being updated, and the appropriate `SELECT` or `ALL` policies will be applied in addition to the `UPDATE` policies. Thus the user must have access to the row(s) being updated through a `SELECT` or `ALL` policy in addition to being granted permission to update the row(s) via an `UPDATE` or `ALL` policy.

When an `INSERT` command has an auxiliary `ON CONFLICT DO UPDATE` clause, if the `UPDATE` path is taken, the row to be updated is first checked against the `USING` expressions of any `UPDATE` policies, and then the new updated row is checked against the `WITH CHECK` expressions. Note, however, that unlike a standalone `UPDATE` command, if the existing row does not pass the `USING` expressions, an error will be thrown (the `UPDATE` path will *never* be silently avoided).

DELETE

Using `DELETE` for a policy means that it will apply to `DELETE` commands. Only rows that pass this policy will be seen by a `DELETE` command. There can be rows that are visible through a `SELECT` that are not available for deletion, if they do not pass the `USING` expression for the `DELETE` policy.

In most cases a `DELETE` command also needs to read data from columns in the relation that it is deleting from (e.g., in a `WHERE` clause or a `RETURNING` clause). In this case, `SELECT` rights are also required on the relation, and the appropriate `SELECT` or `ALL` policies will be applied in addition to the `DELETE` policies. Thus the user must have access to the row(s) being deleted through a `SELECT` or `ALL` policy in addition to being granted permission to delete the row(s) via a `DELETE` or `ALL` policy.

A `DELETE` policy cannot have a `WITH CHECK` expression, as it only applies in cases where records are being deleted from the relation, so that there is no new row to check.

Table 311. Policies Applied by Command Type

Command	SELECT/ALL policy	INSERT/ALL policy	UPDATE/ALL policy		DELETE/ALL policy
	USING expression	WITH CHECK expression	USING expression	WITH CHECK expression	USING expression
SELECT	Existing row	—	—	—	—
SELECT FOR UPDATE/SHARE	Existing row	—	Existing row	—	—
INSERT / MERGE ... THEN INSERT	—	New row	—	—	—
INSERT ... RETURNING	New row ^a	New row	—	—	—
UPDATE / MERGE ... THEN UPDATE	Existing & new rows ^a	—	Existing row	New row	—
DELETE	Existing row ^a	—	—	—	Existing row
ON CONFLICT DO UPDATE	Existing & new rows	—	Existing row	New row	—

^a If read access is required to the existing or new row (for example, a WHERE or RETURNING clause that refers to columns from the relation).

Application of Multiple Policies

When multiple policies of different command types apply to the same command (for example, `SELECT` and `UPDATE` policies applied to an `UPDATE` command), then the user must have both types of permissions (for example, permission to select rows from the relation as well as permission to update them). Thus the expressions for one type of policy are combined with the expressions for the other type of policy using the `AND` operator.

When multiple policies of the same command type apply to the same command, then there must be at least one `PERMISSIVE` policy granting access to the relation, and all of the `RESTRICTIVE` policies must pass. Thus all the `PERMISSIVE` policy expressions are combined using `OR`, all the `RESTRICTIVE` policy expressions are combined using `AND`, and the results are combined using `AND`. If there are no `PERMISSIVE` policies, then access is denied.

Note that, for the purposes of combining multiple policies, `ALL` policies are treated as having the same type as whichever other type of policy is being applied.

For example, in an `UPDATE` command requiring both `SELECT` and `UPDATE` permissions, if there are multiple applicable policies of each type, they will be combined as follows:

```
expression from RESTRICTIVE SELECT/ALL policy 1
AND
expression from RESTRICTIVE SELECT/ALL policy 2
AND
...
AND
(
  expression from PERMISSIVE SELECT/ALL policy 1
  OR
  expression from PERMISSIVE SELECT/ALL policy 2
  OR
  ...
)
AND
```

```
expression from RESTRICTIVE UPDATE/ALL policy 1
AND
expression from RESTRICTIVE UPDATE/ALL policy 2
AND
...
AND
(
    expression from PERMISSIVE UPDATE/ALL policy 1
    OR
    expression from PERMISSIVE UPDATE/ALL policy 2
    OR
    ...
)
```

Notes

You must be the owner of a table to create or change policies for it.

While policies will be applied for explicit queries against tables in the database, they are not applied when the system is performing internal referential integrity checks or validating constraints. This means there are indirect ways to determine that a given value exists. An example of this is attempting to insert a duplicate value into a column that is a primary key or has a unique constraint. If the insert fails then the user can infer that the value already exists. (This example assumes that the user is permitted by policy to insert records which they are not allowed to see.) Another example is where a user is allowed to insert into a table which references another, otherwise hidden table. Existence can be determined by the user inserting values into the referencing table, where success would indicate that the value exists in the referenced table. These issues can be addressed by carefully crafting policies to prevent users from being able to insert, delete, or update records at all which might possibly indicate a value they are not otherwise able to see, or by using generated values (e.g., surrogate keys) instead of keys with external meanings.

Generally, the system will enforce filter conditions imposed using security policies prior to qualifications that appear in user queries, in order to prevent inadvertent exposure of the protected data to user-defined functions which might not be trustworthy. However, functions and operators marked by the system (or the system administrator) as `LEAKPROOF` may be evaluated before policy expressions, as they are assumed to be trustworthy.

Since policy expressions are added to the user's query directly, they will be run with the rights of the user running the overall query. Therefore, users who are using a given policy must be able to access any tables or functions referenced in the expression or they will simply receive a permission denied error when attempting to query the table that has row-level security enabled. This does not change how views work, however. As with normal queries and views, permission checks and policies for the tables which are referenced by a view will use the view owner's rights and any policies which apply to the view owner, except if the view is defined using the `security_invoker` option (see [CREATE VIEW](#)).

No separate policy exists for `MERGE`. Instead, the policies defined for `SELECT`, `INSERT`, `UPDATE`, and `DELETE` are applied while executing `MERGE`, depending on the actions that are performed.

Additional discussion and practical examples can be found in [Section 5.9](#).

Compatibility

`CREATE POLICY` is a Postgres Pro extension.

See Also

[ALTER POLICY](#), [DROP POLICY](#), [ALTER TABLE](#)

CREATE PROCEDURE

CREATE PROCEDURE — define a new procedure

Synopsis

```
CREATE [ OR REPLACE ] PROCEDURE
    name ( [ [ argmode ] [ argname ] argtype [ { DEFAULT | = } default_expr ]
    [, ...] ] )
    { LANGUAGE lang_name
    | TRANSFORM { FOR TYPE type_name } [, ... ]
    | [ EXTERNAL ] SECURITY INVOKER | [ EXTERNAL ] SECURITY DEFINER
    | SET configuration_parameter { TO value | = value | FROM CURRENT }
    | AS 'definition'
    | AS 'obj_file', 'link_symbol'
    | sql_body
    } ...
```

Description

CREATE PROCEDURE defines a new procedure. CREATE OR REPLACE PROCEDURE will either create a new procedure, or replace an existing definition. To be able to define a procedure, the user must have the USAGE privilege on the language.

If a schema name is included, then the procedure is created in the specified schema. Otherwise it is created in the current schema. The name of the new procedure must not match any existing procedure or function with the same input argument types in the same schema. However, procedures and functions of different argument types can share a name (this is called *overloading*).

To replace the current definition of an existing procedure, use CREATE OR REPLACE PROCEDURE. It is not possible to change the name or argument types of a procedure this way (if you tried, you would actually be creating a new, distinct procedure).

When CREATE OR REPLACE PROCEDURE is used to replace an existing procedure, the ownership and permissions of the procedure do not change. All other procedure properties are assigned the values specified or implied in the command. You must own the procedure to replace it (this includes being a member of the owning role).

The user that creates the procedure becomes the owner of the procedure.

To be able to create a procedure, you must have USAGE privilege on the argument types.

Refer to [Section 39.4](#) for further information on writing procedures.

Parameters

name

The name (optionally schema-qualified) of the procedure to create.

argmode

The mode of an argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN.

argname

The name of an argument.

argtype

The data type(s) of the procedure's arguments (optionally schema-qualified), if any. The argument types can be base, composite, or domain types, or can reference the type of a table column.

Depending on the implementation language it might also be allowed to specify “pseudo-types” such as `cstring`. Pseudo-types indicate that the actual argument type is either incompletely specified, or outside the set of ordinary SQL data types.

The type of a column is referenced by writing `table_name.column_name%TYPE`. Using this feature can sometimes help make a procedure independent of changes to the definition of a table.

default_expr

An expression to be used as default value if the parameter is not specified. The expression has to be coercible to the argument type of the parameter. All input parameters following a parameter with a default value must have default values as well.

lang_name

The name of the language that the procedure is implemented in. It can be `sql`, `c`, `internal`, or the name of a user-defined procedural language, e.g., `plpgsql`. The default is `sql` if `sql_body` is specified. Enclosing the name in single quotes is deprecated and requires matching case.

`TRANSFORM { FOR TYPE type_name } [, ...] }`

Lists which transforms a call to the procedure should apply. Transforms convert between SQL types and language-specific data types; see [CREATE TRANSFORM](#). Procedural language implementations usually have hardcoded knowledge of the built-in types, so those don't need to be listed here. If a procedural language implementation does not know how to handle a type and no transform is supplied, it will fall back to a default behavior for converting data types, but this depends on the implementation.

`[EXTERNAL] SECURITY INVOKER`

`[EXTERNAL] SECURITY DEFINER`

`SECURITY INVOKER` indicates that the procedure is to be executed with the privileges of the user that calls it. That is the default. `SECURITY DEFINER` specifies that the procedure is to be executed with the privileges of the user that owns it.

The key word `EXTERNAL` is allowed for SQL conformance, but it is optional since, unlike in SQL, this feature applies to all procedures not only external ones.

A `SECURITY DEFINER` procedure cannot execute transaction control statements (for example, `COMMIT` and `ROLLBACK`, depending on the language).

configuration_parameter value

The `SET` clause causes the specified configuration parameter to be set to the specified value when the procedure is entered, and then restored to its prior value when the procedure exits. `SET FROM CURRENT` saves the value of the parameter that is current when `CREATE PROCEDURE` is executed as the value to be applied when the procedure is entered.

If a `SET` clause is attached to a procedure, then the effects of a `SET LOCAL` command executed inside the procedure for the same variable are restricted to the procedure: the configuration parameter's prior value is still restored at procedure exit. However, an ordinary `SET` command (without `LOCAL`) overrides the `SET` clause, much as it would do for a previous `SET LOCAL` command: the effects of such a command will persist after procedure exit, unless the current transaction is rolled back.

If a `SET` clause is attached to a procedure, then that procedure cannot execute transaction control statements (for example, `COMMIT` and `ROLLBACK`, depending on the language).

See [SET](#) and [Chapter 19](#) for more information about allowed parameter names and values.

definition

A string constant defining the procedure; the meaning depends on the language. It can be an internal procedure name, the path to an object file, an SQL command, or text in a procedural language.

It is often helpful to use dollar quoting (see [Section 4.1.2.4](#)) to write the procedure definition string, rather than the normal single quote syntax. Without dollar quoting, any single quotes or backslashes in the procedure definition must be escaped by doubling them.

obj_file, link_symbol

This form of the `AS` clause is used for dynamically loadable C language procedures when the procedure name in the C language source code is not the same as the name of the SQL procedure. The string *obj_file* is the name of the shared library file containing the compiled C procedure, and is interpreted as for the `LOAD` command. The string *link_symbol* is the procedure's link symbol, that is, the name of the procedure in the C language source code. If the link symbol is omitted, it is assumed to be the same as the name of the SQL procedure being defined.

When repeated `CREATE PROCEDURE` calls refer to the same object file, the file is only loaded once per session. To unload and reload the file (perhaps during development), start a new session.

sql_body

The body of a `LANGUAGE SQL` procedure. This should be a block

```
BEGIN ATOMIC
  statement;
  statement;
  ...
  statement;
END
```

This is similar to writing the text of the procedure body as a string constant (see *definition* above), but there are some differences: This form only works for `LANGUAGE SQL`, the string constant form works for all languages. This form is parsed at procedure definition time, the string constant form is parsed at execution time; therefore this form cannot support polymorphic argument types and other constructs that are not resolvable at procedure definition time. This form tracks dependencies between the procedure and objects used in the procedure body, so `DROP ... CASCADE` will work correctly, whereas the form using string literals may leave dangling procedures. Finally, this form is more compatible with the SQL standard and other SQL implementations.

Notes

See [CREATE FUNCTION](#) for more details on function creation that also apply to procedures.

Use [CALL](#) to execute a procedure.

Examples

```
CREATE PROCEDURE insert_data(a integer, b integer)
LANGUAGE SQL
AS $$
INSERT INTO tbl VALUES (a);
INSERT INTO tbl VALUES (b);
$$;
```

or

```
CREATE PROCEDURE insert_data(a integer, b integer)
LANGUAGE SQL
BEGIN ATOMIC
  INSERT INTO tbl VALUES (a);
  INSERT INTO tbl VALUES (b);
END;
```

and call like this:

```
CALL insert_data(1, 2);
```

Compatibility

A `CREATE PROCEDURE` command is defined in the SQL standard. The Postgres Pro implementation can be used in a compatible way but has many extensions. For details see also [CREATE FUNCTION](#).

See Also

[ALTER PROCEDURE](#), [DROP PROCEDURE](#), [CALL](#), [CREATE FUNCTION](#)

CREATE PROFILE

CREATE PROFILE — define a new profile

Synopsis

```
CREATE PROFILE [ IF NOT EXISTS ] name [ LIMIT parameter value [ ... ] ]
```

where *parameter* can be:

```
    FAILED_LOGIN_ATTEMPTS
| PASSWORD_REUSE_TIME
| PASSWORD_REUSE_MAX
| PASSWORD_LIFE_TIME
| PASSWORD_GRACE_TIME
| USER_INACTIVE_TIME
| FAILED_AUTH_KEEP_TIME
| PASSWORD_MIN_UNIQUE_CHARS
| PASSWORD_MIN_LEN
| PASSWORD_REQUIRE_COMPLEX
```

```
CREATE PROFILE [ IF NOT EXISTS ] name FROM existing_profile
```

Description

The `CREATE PROFILE` command adds a new profile for the Postgres Pro database cluster. You must be a database superuser or have the privileges of the `pg_manage_profiles` role to use this command.

A *profile* defines a set of parameters that restrict database usage. In particular, Postgres Pro profiles enforce password management policy for the roles assigned to them. Profiles are defined at the database cluster level, so they apply to all databases in the cluster.

By default, all parameter values of a new profile are set to `DEFAULT`, which inherits the actual values from the built-in `default` profile. Initially, the `default` profile provides no usage restrictions, but it can be changed by the `ALTER PROFILE` command. The `UNLIMITED` value indicates that no restrictions apply to a particular parameter.

Once a profile is assigned to a role, all restrictions of this profile apply to this role. All newly created roles have the `default` profile, unless you explicitly assign a different profile to this role.

Parameters

name

The name of the new profile.

`FAILED_LOGIN_ATTEMPTS` *value*

Specifies the number of failed login attempts before the role is locked. A superuser can unlock the locked role by running the `ALTER ROLE` command with the `ACCOUNT UNLOCK` clause.

Note that there can be several actual login attempts made behind each user-perceived login attempt. For example, when the user tries to log in with SSL enabled, libpq-based clients by default also make a non-SSL connection attempt if an SSL connection fails.

Possible values are integers greater than 0, `DEFAULT`, or `UNLIMITED`.

PASSWORD_REUSE_TIME *value*

Specifies for how long an old password cannot be reused. Measured in days. Possible values are real numbers greater than or equal to 0, interval values, **DEFAULT**, or **UNLIMITED**.

Set this parameter together with **PASSWORD_REUSE_MAX** as its effect depends on the combination. If both these parameters are set to **UNLIMITED**, there are no restrictions on password reuse. If only one of them is set to **UNLIMITED**, password reuse is always forbidden.

PASSWORD_REUSE_MAX *value*

Specifies the number of password changes required before the current password can be reused. Possible values are integers greater than or equal to 0, **DEFAULT**, or **UNLIMITED**.

Set this parameter together with **PASSWORD_REUSE_TIME** as its effect depends on the combination. If both these parameters are set to **UNLIMITED**, there are no restrictions on password reuse. If only one of them is set to **UNLIMITED**, password reuse is always forbidden.

PASSWORD_LIFE_TIME *value*

Specifies for how long the password can be used for authentication on the primary server. Measured in days. Possible values are real numbers, interval values, **DEFAULT**, or **UNLIMITED**. The resulting value in seconds must be greater than 0. Once the password expires, all further connections of the corresponding role are rejected. The role can be unlocked with the [ALTER ROLE](#) command.

However, if the [LDAP authentication](#) is configured, this role still can connect to standby servers, unless it is locked in LDAP too.

If **PASSWORD_GRACE_TIME** parameter is also set, the specified grace period is added to the password lifetime. During the grace period, the role is prompted to change the password while still allowed to log in.

PASSWORD_GRACE_TIME *value*

Specifies for how long a warning is raised that the password is going to expire while login is still allowed, measured in days. You can set the password life time in the **VALID UNTIL** attribute of the role or in the **PASSWORD_LIFE_TIME** parameter of the profile.

Possible values are real numbers greater than or equal to 0, interval values, **DEFAULT**, or **UNLIMITED**. If the **PASSWORD_GRACE_TIME** parameter is set to **UNLIMITED**, the password life time effectively becomes unlimited.

USER_INACTIVE_TIME *value*

Specifies for how long the user can be inactive since the last login before the role is locked. Measured in days. A superuser can unlock the locked role by running the [ALTER ROLE](#) command with the **ACCOUNT UNLOCK** clause. Possible values are real numbers, interval values, **DEFAULT**, or **UNLIMITED**. The resulting value in seconds must be greater than 0.

FAILED_AUTH_KEEP_TIME *value*

Specifies for how long the information on the user's first authentication failure is kept. Measured in days. Possible values are real numbers, interval values, **DEFAULT**, or **UNLIMITED**. The resulting value in seconds must be greater than 0. When the user attempts to log in after this time interval expires, the failed login attempts counter (see **FAILED_LOGIN_ATTEMPTS** parameter) is reset and the user is unlocked if they were locked previously due to authentication failures.

PASSWORD_MIN_UNIQUE_CHARS *value*

Specifies minimum number of unique characters for a password. Possible values are integers greater than 0, **DEFAULT**, or **UNLIMITED**.

`PASSWORD_MIN_LEN` *value*

Specifies minimum number of characters for a password. Possible values are positive integers, `DEFAULT`, or `UNLIMITED`.

`PASSWORD_REQUIRE_COMPLEX` [*value*]

Specifies whether password complexity is checked. If this check is enabled, a password must meet the following requirements:

- Password contains at least one character from three of the following groups: lowercase letters, uppercase letters, digits, and special characters
- Password doesn't contain the user name

Possible values are booleans or `DEFAULT`. If the parameter is used without a value, it is set as `true`.

`IF NOT EXISTS`

Do not throw an error if a profile with the same name already exists.

existing_profile

The name of an existing profile to copy. The new profile will have the same properties as the existing one, but it will be an independent object.

Important

All letters of languages without case distinction (Hindi, Chinese, etc.) in the UTF-8 encoding are considered lowercase.

Notes

Use [ALTER PROFILE](#) to change the parameter values of a profile, and [DROP PROFILE](#) to remove a profile. All the parameters specified by `CREATE PROFILE` can be modified later by running the `ALTER PROFILE` command.

Warning

`PASSWORD_REUSE_TIME` and `PASSWORD_REUSE_MAX` parameters might not work if the password is passed in an encrypted form during the password change. The `\password` command of `psql` encrypts the password (see [the section called “Meta-Commands”](#) for details). If the password is MD5 encrypted, it can be checked for equivalence, and `PASSWORD_REUSE_TIME` or `PASSWORD_REUSE_MAX` parameters can be applied. But there is no way to check for equivalence for SCRAM-SHA-256 encrypted passwords.

Warning

`PASSWORD_MIN_UNIQUE_CHARS`, `PASSWORD_MIN_LEN`, and `PASSWORD_REQUIRE_COMPLEX` parameters don't work if the password is passed in an encrypted form during the password change.

Examples

Create a profile `admin_profile`:

```
CREATE PROFILE admin_profile
    LIMIT PASSWORD_REUSE_MAX 10
    PASSWORD_REUSE_TIME 30;
```

Create a role having `admin_profile`:

```
CREATE ROLE admin WITH PROFILE admin_profile;
```

Create a profile from an existing profile:

```
CREATE PROFILE administrator FROM admin_profile;
```

This can be convenient to be able to use an existing profile as a template for a new one.

See Also

[ALTER PROFILE](#), [DROP PROFILE](#), [CREATE ROLE](#)

CREATE PUBLICATION

CREATE PUBLICATION — define a new publication

Synopsis

```
CREATE PUBLICATION name
  [ FOR ALL TABLES
    | FOR publication_object [, ... ] ]
  [ WITH ( publication_parameter [= value] [, ... ] ) ]
```

where *publication_object* is one of:

```
TABLE [ ONLY ] table_name [ * ] [ ( column_name [, ... ] ) ] [ WHERE ( expression
) ] [, ... ]
TABLES IN SCHEMA { schema_name | CURRENT_SCHEMA } [, ... ]
```

Description

CREATE PUBLICATION adds a new publication into the current database. The publication name must be distinct from the name of any existing publication in the current database.

A publication is essentially a group of tables whose data changes are intended to be replicated through logical replication. See [Section 29.1](#) for details about how publications fit into the logical replication setup.

Parameters

name

The name of the new publication.

FOR TABLE

Specifies a list of tables to add to the publication. If ONLY is specified before the table name, only that table is added to the publication. If ONLY is not specified, the table and all its descendant tables (if any) are added. Optionally, * can be specified after the table name to explicitly indicate that descendant tables are included. This does not apply to a partitioned table, however. The partitions of a partitioned table are always implicitly considered part of the publication, so they are never explicitly added to the publication.

If the optional WHERE clause is specified, it defines a *row filter* expression. Rows for which the *expression* evaluates to false or null will not be published. Note that parentheses are required around the expression. It has no effect on TRUNCATE commands.

When a column list is specified, only the named columns are replicated. If no column list is specified, all columns of the table are replicated through this publication, including any columns added later. It has no effect on TRUNCATE commands. See [Section 29.5](#) for details about column lists.

Only persistent base tables and partitioned tables can be part of a publication. Temporary tables, unlogged tables, foreign tables, materialized views, and regular views cannot be part of a publication.

Specifying a column list when the publication also publishes FOR TABLES IN SCHEMA is not supported.

When a partitioned table is added to a publication, all of its existing and future partitions are implicitly considered to be part of the publication. So, even operations that are performed directly on a partition are also published via publications that its ancestors are part of.

FOR ALL TABLES

Marks the publication as one that replicates changes for all tables in the database, including tables created in the future.

FOR TABLES IN SCHEMA

Marks the publication as one that replicates changes for all tables in the specified list of schemas, including tables created in the future.

Specifying a schema when the publication also publishes a table with a column list is not supported.

Only persistent base tables and partitioned tables present in the schema will be included as part of the publication. Temporary tables, unlogged tables, foreign tables, materialized views, and regular views from the schema will not be part of the publication.

When a partitioned table is published via schema level publication, all of its existing and future partitions are implicitly considered to be part of the publication, regardless of whether they are from the publication schema or not. So, even operations that are performed directly on a partition are also published via publications that its ancestors are part of.

WITH (*publication_parameter* [= *value*] [, ...])

This clause specifies optional parameters for a publication. The following parameters are supported:

publish (*string*)

This parameter determines which DML operations will be published by the new publication to the subscribers. The value is comma-separated list of operations. The allowed operations are `insert`, `update`, `delete`, and `truncate`. The default is to publish all actions, and so the default value for this option is `'insert, update, delete, truncate'`.

This parameter only affects DML operations. In particular, the initial data synchronization (see [Section 29.8.1](#)) for logical replication does not take this parameter into account when copying existing table data.

publish_via_partition_root (*boolean*)

This parameter determines whether changes in a partitioned table (or on its partitions) contained in the publication will be published using the identity and schema of the partitioned table rather than that of the individual partitions that are actually changed; the latter is the default. Enabling this allows the changes to be replicated into a non-partitioned table or a partitioned table consisting of a different set of partitions.

There can be a case where a subscription combines multiple publications. If a partitioned table is published by any subscribed publications which set `publish_via_partition_root = true`, changes on this partitioned table (or on its partitions) will be published using the identity and schema of this partitioned table rather than that of the individual partitions.

This parameter also affects how row filters and column lists are chosen for partitions; see below for details.

If this is enabled, `TRUNCATE` operations performed directly on partitions are not replicated.

When specifying a parameter of type `boolean`, the `= value` part can be omitted, which is equivalent to specifying `TRUE`.

Notes

If `FOR TABLE`, `FOR ALL TABLES` or `FOR TABLES IN SCHEMA` are not specified, then the publication starts out with an empty set of tables. That is useful if tables or schemas are to be added later.

The creation of a publication does not start replication. It only defines a grouping and filtering logic for future subscribers.

To create a publication, the invoking user must have the `CREATE` privilege for the current database. (Of course, superusers bypass this check.)

To add a table to a publication, the invoking user must have ownership rights on the table. The `FOR ALL TABLES` and `FOR TABLES IN SCHEMA` clauses require the invoking user to be a superuser.

The tables added to a publication that publishes `UPDATE` and/or `DELETE` operations must have `REPLICA IDENTITY` defined. Otherwise those operations will be disallowed on those tables.

Any column list must include the `REPLICA IDENTITY` columns in order for `UPDATE` or `DELETE` operations to be published. There are no column list restrictions if the publication publishes only `INSERT` operations.

A row filter expression (i.e., the `WHERE` clause) must contain only columns that are covered by the `REPLICA IDENTITY`, in order for `UPDATE` and `DELETE` operations to be published. For publication of `INSERT` operations, any column may be used in the `WHERE` expression. The row filter allows simple expressions that don't have user-defined functions, user-defined operators, user-defined types, user-defined collations, non-immutable built-in functions, or references to system columns.

The row filter on a table becomes redundant if `FOR TABLES IN SCHEMA` is specified and the table belongs to the referred schema.

For published partitioned tables, the row filter for each partition is taken from the published partitioned table if the publication parameter `publish_via_partition_root` is true, or from the partition itself if it is false (the default). See [Section 29.4](#) for details about row filters. Similarly, for published partitioned tables, the column list for each partition is taken from the published partitioned table if the publication parameter `publish_via_partition_root` is true, or from the partition itself if it is false.

For an `INSERT ... ON CONFLICT` command, the publication will publish the operation that results from the command. Depending on the outcome, it may be published as either `INSERT` or `UPDATE`, or it may not be published at all.

For a `MERGE` command, the publication will publish an `INSERT`, `UPDATE`, or `DELETE` for each row inserted, updated, or deleted.

ATTACHing a table into a partition tree whose root is published using a publication with `publish_via_partition_root` set to true does not result in the table's existing contents being replicated.

`COPY ... FROM` commands are published as `INSERT` operations.

DDL operations are not published.

The `WHERE` clause expression is executed with the role used for the replication connection.

Examples

Create a publication that publishes all changes in two tables:

```
CREATE PUBLICATION mypublication FOR TABLE users, departments;
```

Create a publication that publishes all changes from active departments:

```
CREATE PUBLICATION active_departments FOR TABLE departments WHERE (active IS TRUE);
```

Create a publication that publishes all changes in all tables:

```
CREATE PUBLICATION alltables FOR ALL TABLES;
```

Create a publication that only publishes `INSERT` operations in one table:

```
CREATE PUBLICATION insert_only FOR TABLE mydata
WITH (publish = 'insert');
```

Create a publication that publishes all changes for tables `users`, `departments` and all changes for all the tables present in the schema `production`:

```
CREATE PUBLICATION production_publication FOR TABLE users, departments, TABLES IN
  SCHEMA production;
```

Create a publication that publishes all changes for all the tables present in the schemas `marketing` and `sales`:

```
CREATE PUBLICATION sales_publication FOR TABLES IN SCHEMA marketing, sales;
```

Create a publication that publishes all changes for table `users`, but replicates only columns `user_id` and `firstname`:

```
CREATE PUBLICATION users_filtered FOR TABLE users (user_id, firstname);
```

Compatibility

`CREATE PUBLICATION` is a Postgres Pro extension.

See Also

[ALTER PUBLICATION](#), [DROP PUBLICATION](#), [CREATE SUBSCRIPTION](#), [ALTER SUBSCRIPTION](#)

CREATE ROLE

CREATE ROLE — define a new database role

Synopsis

```
CREATE ROLE name [ [ WITH ] option [ ... ] ]
```

where *option* can be:

```
SUPERUSER | NOSUPERUSER
| CREATEDB | NOCREATEDB
| CREATEROLE | NOCREATEROLE
| INHERIT | NOINHERIT
| LOGIN | NOLOGIN
| REPLICATION | NOREPLICATION
| BYPASSRLS | NOBYPASSRLS
| CONNECTION LIMIT connlimit
| [ ENCRYPTED ] PASSWORD 'password' | PASSWORD NULL
| VALID UNTIL 'timestamp'
| IN ROLE role_name [, ...]
| ROLE role_name [, ...]
| ADMIN role_name [, ...]
| SYSID uid
| PROFILE profile_name
```

Description

CREATE ROLE adds a new role to a Postgres Pro database cluster. A role is an entity that can own database objects and have database privileges; a role can be considered a “user”, a “group”, or both depending on how it is used. Refer to [Chapter 21](#) and [Chapter 20](#) for information about managing users and authentication. You must have CREATEROLE privilege or be a database superuser to use this command.

Note that roles are defined at the database cluster level, and so are valid in all databases in the cluster.

During role creation it is possible to immediately assign the newly created role to be a member of an existing role, and also assign existing roles to be members of the newly created role. The rules for which initial role membership options are enabled are described below in the IN ROLE, ROLE, and ADMIN clauses. The GRANT command has fine-grained option control during membership creation, and the ability to modify these options after the new role is created.

Parameters

name

The name of the new role.

SUPERUSER

NOSUPERUSER

These clauses determine whether the new role is a “superuser”, who can override all access restrictions within the database. Superuser status is dangerous and should be used only when really needed. You must yourself be a superuser to create a new superuser. If not specified, NOSUPERUSER is the default.

CREATEDB

NOCREATEDB

These clauses define a role's ability to create databases. If CREATEDB is specified, the role being defined will be allowed to create new databases. Specifying NOCREATEDB will deny a role the ability

to create databases. If not specified, `NOCREATEDB` is the default. Only superuser roles or roles with `CREATEDB` can specify `CREATEDB`.

`CREATEROLE`
`NOCREATEROLE`

These clauses determine whether a role will be permitted to create, alter, drop, comment on, and change the security label for other roles. See [role creation](#) for more details about what capabilities are conferred by this privilege. If not specified, `NOCREATEROLE` is the default.

`INHERIT`
`NOINHERIT`

This affects the membership inheritance status when this role is added as a member of another role, both in this and future commands. Specifically, it controls the inheritance status of memberships added with this command using the `IN ROLE` clause, and in later commands using the `ROLE` clause. It is also used as the default inheritance status when adding this role as a member using the `GRANT` command. If not specified, `INHERIT` is the default.

In Postgres Pro versions before 16, inheritance was a role-level attribute that controlled all runtime membership checks for that role.

`LOGIN`
`NOLOGIN`

These clauses determine whether a role is allowed to log in; that is, whether the role can be given as the initial session authorization name during client connection. A role having the `LOGIN` attribute can be thought of as a user. Roles without this attribute are useful for managing database privileges, but are not users in the usual sense of the word. If not specified, `NOLOGIN` is the default, except when `CREATE ROLE` is invoked through its alternative spelling `CREATE USER`.

`REPLICATION`
`NOREPLICATION`

These clauses determine whether a role is a replication role. A role must have this attribute (or be a superuser) in order to be able to connect to the server in replication mode (physical or logical replication) and in order to be able to create or drop replication slots. A role having the `REPLICATION` attribute is a very highly privileged role, and should only be used on roles actually used for replication. If not specified, `NOREPLICATION` is the default. Only superuser roles or roles with `REPLICATION` can specify `REPLICATION`.

`BYPASSRLS`
`NOBYPASSRLS`

These clauses determine whether a role bypasses every row-level security (RLS) policy. `NOBYPASSRLS` is the default. Only superuser roles or roles with `BYPASSRLS` can specify `BYPASSRLS`.

Note that `pg_dump` will set `row_security` to `OFF` by default, to ensure all contents of a table are dumped out. If the user running `pg_dump` does not have appropriate permissions, an error will be returned. However, superusers and the owner of the table being dumped always bypass RLS.

`CONNECTION LIMIT` *conlimit*

If role can log in, this specifies how many concurrent connections the role can make. `-1` (the default) means no limit. Note that only normal connections are counted towards this limit. Neither prepared transactions nor background worker connections are counted towards this limit.

`[ENCRYPTED] PASSWORD 'password'`
`PASSWORD NULL`

Sets the role's password. (A password is only of use for roles having the `LOGIN` attribute, but you can nonetheless define one for roles without it.) If you do not plan to use password authentication

you can omit this option. If no password is specified, the password will be set to null and password authentication will always fail for that user. A null password can optionally be written explicitly as `PASSWORD NULL`.

Note

Specifying an empty string will also set the password to null, but that was not the case before Postgres Pro version 10. In earlier versions, an empty string could be used, or not, depending on the authentication method and the exact version, and libpq would refuse to use it in any case. To avoid the ambiguity, specifying an empty string should be avoided.

The password is always stored encrypted in the system catalogs. The `ENCRYPTED` keyword has no effect, but is accepted for backwards compatibility. The method of encryption is determined by the configuration parameter `password_encryption`. If the presented password string is already in MD5-encrypted or SCRAM-encrypted format, then it is stored as-is regardless of `password_encryption` (since the system cannot decrypt the specified encrypted password string, to encrypt it in a different format). This allows reloading of encrypted passwords during dump/restore.

`VALID UNTIL 'timestamp'`

The `VALID UNTIL` clause sets a date and time after which the role's password is no longer valid. If this clause is omitted the password will be valid for all time.

`IN ROLE role_name`

The `IN ROLE` clause causes the new role to be automatically added as a member of the specified existing roles. The new membership will have the `SET` option enabled and the `ADMIN` option disabled. The `INHERIT` option will be enabled unless the `NOINHERIT` option is specified.

`ROLE role_name`

The `ROLE` clause causes one or more specified existing roles to be automatically added as members, with the `SET` option enabled. This in effect makes the new role a “group”. Roles named in this clause with the role-level `INHERIT` attribute will have the `INHERIT` option enabled in the new membership. New memberships will have the `ADMIN` option disabled.

`ADMIN role_name`

The `ADMIN` clause has the same effect as `ROLE`, but the named roles are added as members of the new role with `ADMIN` enabled, giving them the right to grant membership in the new role to others.

`SYSID uid`

The `SYSID` clause is ignored, but is accepted for backwards compatibility.

`PROFILE profile_name`

Sets the role's profile. If this clause is omitted, the `default` profile is assigned to a role (see [CREATE PROFILE](#) for details).

Notes

Use [ALTER ROLE](#) to change the attributes of a role, and [DROP ROLE](#) to remove a role. All the attributes specified by `CREATE ROLE` can be modified by later `ALTER ROLE` commands.

The preferred way to add and remove members of roles that are being used as groups is to use [GRANT](#) and [REVOKE](#).

The `VALID UNTIL` clause defines an expiration time for a password only, not for the role per se. In particular, the expiration time is not enforced when logging in using a non-password-based authentication method.

The role attributes defined here are non-inheritable, i.e., being a member of a role with, e.g., `CREATEDB` will not allow the member to create new databases even if the membership grant has the `INHERIT` option. Of course, if the membership grant has the `SET` option the member role would be able to `SET ROLE` to the `createdb` role and then create a new database.

The membership grants created by the `IN ROLE`, `ROLE`, and `ADMIN` clauses have the role executing this command as the grantor.

The `INHERIT` attribute is the default for reasons of backwards compatibility: in prior releases of Postgres Pro, users always had access to all privileges of groups they were members of. However, `NOINHERIT` provides a closer match to the semantics specified in the SQL standard.

Postgres Pro includes a program `createuser` that has the same functionality as `CREATE ROLE` (in fact, it calls this command) but can be run from the command shell.

The `CONNECTION LIMIT` option is only enforced approximately; if two new sessions start at about the same time when just one connection “slot” remains for the role, it is possible that both will fail. Also, the limit is never enforced for superusers.

Caution must be exercised when specifying an unencrypted password with this command. The password will be transmitted to the server in cleartext, and it might also be logged in the client's command history or the server log. The command `createuser`, however, transmits the password encrypted. Also, `psql` contains a command `\password` that can be used to safely change the password later.

Examples

Create a role that can log in, but don't give it a password:

```
CREATE ROLE jonathan LOGIN;
```

Create a role with a password:

```
CREATE USER davide WITH PASSWORD 'jw8s0F4';
```

(`CREATE USER` is the same as `CREATE ROLE` except that it implies `LOGIN`.)

Create a role with a profile:

```
CREATE ROLE davide WITH PROFILE admin_profile;
```

Create a role with a password that is valid until the end of 2004. After one second has ticked in 2005, the password is no longer valid.

```
CREATE ROLE miriam WITH LOGIN PASSWORD 'jw8s0F4' VALID UNTIL '2005-01-01';
```

Create a role that can create databases and manage roles:

```
CREATE ROLE admin WITH CREATEDB CREATEROLE;
```

Compatibility

The `CREATE ROLE` statement is in the SQL standard, but the standard only requires the syntax

```
CREATE ROLE name [ WITH ADMIN role_name ]
```

Multiple initial administrators, and all the other options of `CREATE ROLE`, are Postgres Pro extensions.

The SQL standard defines the concepts of users and roles, but it regards them as distinct concepts and leaves all commands defining users to be specified by each database implementation. In Postgres Pro we have chosen to unify users and roles into a single kind of entity. Roles therefore have many more optional attributes than they do in the standard.

The behavior specified by the SQL standard is most closely approximated creating SQL-standard users as Postgres Pro roles with the `NOINHERIT` option, and SQL-standard roles as Postgres Pro roles with the `INHERIT` option.

The `USER` clause has the same behavior as `ROLE` but has been deprecated:

```
USER role_name [, ...]
```

The `IN GROUP` clause has the same behavior as `IN ROLE` but has been deprecated:

```
IN GROUP role_name [, ...]
```

See Also

[SET ROLE](#), [ALTER ROLE](#), [DROP ROLE](#), [GRANT](#), [REVOKE](#), [createuser](#), [createrole_self_grant](#), [CREATE PROFILE](#)

CREATE RULE

CREATE RULE — define a new rewrite rule

Synopsis

```
CREATE [ OR REPLACE ] RULE name AS ON event
  TO table_name [ WHERE condition ]
  DO [ ALSO | INSTEAD ] { NOTHING | command | ( command ; command ... ) }
```

where *event* can be one of:

```
SELECT | INSERT | UPDATE | DELETE
```

Description

CREATE RULE defines a new rule applying to a specified table or view. CREATE OR REPLACE RULE will either create a new rule, or replace an existing rule of the same name for the same table.

The Postgres Pro rule system allows one to define an alternative action to be performed on insertions, updates, or deletions in database tables. Roughly speaking, a rule causes additional commands to be executed when a given command on a given table is executed. Alternatively, an INSTEAD rule can replace a given command by another, or cause a command not to be executed at all. Rules are used to implement SQL views as well. It is important to realize that a rule is really a command transformation mechanism, or command macro. The transformation happens before the execution of the command starts. If you actually want an operation that fires independently for each physical row, you probably want to use a trigger, not a rule. More information about the rules system is in [Chapter 42](#).

Presently, ON SELECT rules can only be attached to views. Such a rule must be named "_RETURN", must be an unconditional INSTEAD rule, and must have an action that consists of a single SELECT command. This command defines the visible contents of the view. (The view itself is basically a dummy table with no storage.) It's best to regard such a rule as an implementation detail. While a view can be redefined via CREATE OR REPLACE RULE "_RETURN" AS ..., it's better style to use CREATE OR REPLACE VIEW.

You can create the illusion of an updatable view by defining ON INSERT, ON UPDATE, and ON DELETE rules (or any subset of those that's sufficient for your purposes) to replace update actions on the view with appropriate updates on other tables. If you want to support INSERT RETURNING and so on, then be sure to put a suitable RETURNING clause into each of these rules.

There is a catch if you try to use conditional rules for complex view updates: there *must* be an unconditional INSTEAD rule for each action you wish to allow on the view. If the rule is conditional, or is not INSTEAD, then the system will still reject attempts to perform the update action, because it thinks it might end up trying to perform the action on the dummy table of the view in some cases. If you want to handle all the useful cases in conditional rules, add an unconditional DO INSTEAD NOTHING rule to ensure that the system understands it will never be called on to update the dummy table. Then make the conditional rules non-INSTEAD; in the cases where they are applied, they add to the default INSTEAD NOTHING action. (This method does not currently work to support RETURNING queries, however.)

Note

A view that is simple enough to be automatically updatable (see [CREATE VIEW](#)) does not require a user-created rule in order to be updatable. While you can create an explicit rule anyway, the automatic update transformation will generally outperform an explicit rule.

Another alternative worth considering is to use INSTEAD OF triggers (see [CREATE TRIGGER](#)) in place of rules.

Parameters

name

The name of a rule to create. This must be distinct from the name of any other rule for the same table. Multiple rules on the same table and same event type are applied in alphabetical name order.

event

The event is one of `SELECT`, `INSERT`, `UPDATE`, or `DELETE`. Note that an `INSERT` containing an `ON CONFLICT` clause cannot be used on tables that have either `INSERT` or `UPDATE` rules. Consider using an updatable view instead.

table_name

The name (optionally schema-qualified) of the table or view the rule applies to.

condition

Any SQL conditional expression (returning `boolean`). The condition expression cannot refer to any tables except `NEW` and `OLD`, and cannot contain aggregate functions.

`INSTEAD`

`INSTEAD` indicates that the commands should be executed *instead of* the original command.

`ALSO`

`ALSO` indicates that the commands should be executed *in addition to* the original command.

If neither `ALSO` nor `INSTEAD` is specified, `ALSO` is the default.

command

The command or commands that make up the rule action. Valid commands are `SELECT`, `INSERT`, `UPDATE`, `DELETE`, or `NOTIFY`.

Within *condition* and *command*, the special table names `NEW` and `OLD` can be used to refer to values in the referenced table. `NEW` is valid in `ON INSERT` and `ON UPDATE` rules to refer to the new row being inserted or updated. `OLD` is valid in `ON UPDATE` and `ON DELETE` rules to refer to the existing row being updated or deleted.

Notes

You must be the owner of a table to create or change rules for it.

In a rule for `INSERT`, `UPDATE`, or `DELETE` on a view, you can add a `RETURNING` clause that emits the view's columns. This clause will be used to compute the outputs if the rule is triggered by an `INSERT RETURNING`, `UPDATE RETURNING`, or `DELETE RETURNING` command respectively. When the rule is triggered by a command without `RETURNING`, the rule's `RETURNING` clause will be ignored. The current implementation allows only unconditional `INSTEAD` rules to contain `RETURNING`; furthermore there can be at most one `RETURNING` clause among all the rules for the same event. (This ensures that there is only one candidate `RETURNING` clause to be used to compute the results.) `RETURNING` queries on the view will be rejected if there is no `RETURNING` clause in any available rule.

It is very important to take care to avoid circular rules. For example, though each of the following two rule definitions are accepted by Postgres Pro, the `SELECT` command would cause Postgres Pro to report an error because of recursive expansion of a rule:

```
CREATE RULE "_RETURN" AS
  ON SELECT TO t1
  DO INSTEAD
    SELECT * FROM t2;
```

```
CREATE RULE "_RETURN" AS
  ON SELECT TO t2
  DO INSTEAD
    SELECT * FROM t1;
```

```
SELECT * FROM t1;
```

Presently, if a rule action contains a `NOTIFY` command, the `NOTIFY` command will be executed unconditionally, that is, the `NOTIFY` will be issued even if there are not any rows that the rule should apply to. For example, in:

```
CREATE RULE notify_me AS ON UPDATE TO mytable DO ALSO NOTIFY mytable;
```

```
UPDATE mytable SET name = 'foo' WHERE id = 42;
```

one `NOTIFY` event will be sent during the `UPDATE`, whether or not there are any rows that match the condition `id = 42`. This is an implementation restriction that might be fixed in future releases.

Compatibility

`CREATE RULE` is a Postgres Pro language extension, as is the entire query rewrite system.

See Also

[ALTER RULE](#), [DROP RULE](#)

CREATE SCHEMA

CREATE SCHEMA — define a new schema

Synopsis

```
CREATE SCHEMA schema_name [ AUTHORIZATION role_specification ] [ schema_element
[ ... ] ]
CREATE SCHEMA AUTHORIZATION role_specification [ schema_element [ ... ] ]
CREATE SCHEMA IF NOT EXISTS schema_name [ AUTHORIZATION role_specification ]
CREATE SCHEMA IF NOT EXISTS AUTHORIZATION role_specification
```

where *role_specification* can be:

```
user_name
| CURRENT_ROLE
| CURRENT_USER
| SESSION_USER
```

Description

CREATE SCHEMA enters a new schema into the current database. The schema name must be distinct from the name of any existing schema in the current database.

A schema is essentially a namespace: it contains named objects (tables, data types, functions, and operators) whose names can duplicate those of other objects existing in other schemas. Named objects are accessed either by “qualifying” their names with the schema name as a prefix, or by setting a search path that includes the desired schema(s). A CREATE command specifying an unqualified object name creates the object in the current schema (the one at the front of the search path, which can be determined with the function `current_schema()`).

Optionally, CREATE SCHEMA can include subcommands to create objects within the new schema. The subcommands are treated essentially the same as separate commands issued after creating the schema, except that if the AUTHORIZATION clause is used, all the created objects will be owned by that user.

Parameters

schema_name

The name of a schema to be created. If this is omitted, the *user_name* is used as the schema name. The name cannot begin with `pg_`, as such names are reserved for system schemas.

user_name

The role name of the user who will own the new schema. If omitted, defaults to the user executing the command. To create a schema owned by another role, you must be able to SET ROLE to that role.

schema_element

An SQL statement defining an object to be created within the schema. Currently, only CREATE TABLE, CREATE VIEW, CREATE INDEX, CREATE SEQUENCE, CREATE TRIGGER and GRANT are accepted as clauses within CREATE SCHEMA. Other kinds of objects may be created in separate commands after the schema is created.

IF NOT EXISTS

Do nothing (except issuing a notice) if a schema with the same name already exists. *schema_element* subcommands cannot be included when this option is used.

Notes

To create a schema, the invoking user must have the `CREATE` privilege for the current database. (Of course, superusers bypass this check.)

Examples

Create a schema:

```
CREATE SCHEMA myschema;
```

Create a schema for user `joe`; the schema will also be named `joe`:

```
CREATE SCHEMA AUTHORIZATION joe;
```

Create a schema named `test` that will be owned by user `joe`, unless there already is a schema named `test`. (It does not matter whether `joe` owns the pre-existing schema.)

```
CREATE SCHEMA IF NOT EXISTS test AUTHORIZATION joe;
```

Create a schema and create a table and view within it:

```
CREATE SCHEMA hollywood
    CREATE TABLE films (title text, release date, awards text[])
    CREATE VIEW winners AS
        SELECT title, release FROM films WHERE awards IS NOT NULL;
```

Notice that the individual subcommands do not end with semicolons.

The following is an equivalent way of accomplishing the same result:

```
CREATE SCHEMA hollywood;
CREATE TABLE hollywood.films (title text, release date, awards text[]);
CREATE VIEW hollywood.winners AS
    SELECT title, release FROM hollywood.films WHERE awards IS NOT NULL;
```

Compatibility

The SQL standard allows a `DEFAULT CHARACTER SET` clause in `CREATE SCHEMA`, as well as more subcommand types than are presently accepted by Postgres Pro.

The SQL standard specifies that the subcommands in `CREATE SCHEMA` can appear in any order. The present Postgres Pro implementation does not handle all cases of forward references in subcommands; it might sometimes be necessary to reorder the subcommands in order to avoid forward references.

According to the SQL standard, the owner of a schema always owns all objects within it. Postgres Pro allows schemas to contain objects owned by users other than the schema owner. This can happen only if the schema owner grants the `CREATE` privilege on their schema to someone else, or a superuser chooses to create objects in it.

The `IF NOT EXISTS` option is a Postgres Pro extension.

See Also

[ALTER SCHEMA](#), [DROP SCHEMA](#)

CREATE SEQUENCE

CREATE SEQUENCE — define a new sequence generator

Synopsis

For sharded objects:

```
CREATE SEQUENCE [ IF NOT EXISTS ] name
  [ AS data_type ]
  [ INCREMENT [ BY ] increment ]
  [ MINVALUE minvalue | NO MINVALUE ] [ MAXVALUE maxvalue | NO MAXVALUE ]

  [ CACHE cache ] [ [ NO CYCLE ]
  [ OWNED BY { table_name.column_name | NONE } ]
WITH ( [ global = true ],
      [ block_size = block_size ]
      )
```

For other object types:

```
CREATE [ { TEMPORARY | TEMP } | UNLOGGED ] SEQUENCE [ IF NOT EXISTS ] name
  [ AS data_type ]
  [ INCREMENT [ BY ] increment ]
  [ MINVALUE minvalue | NO MINVALUE ] [ MAXVALUE maxvalue | NO MAXVALUE ]
  [ START [ WITH ] start ] [ CACHE cache ] [ [ NO ] CYCLE ]
  [ OWNED BY { table_name.column_name | NONE } ]
```

Description

CREATE SEQUENCE creates a new sequence number generator. This involves creating and initializing a new special single-row table with the name *name*. The generator will be owned by the user issuing the command.

Postgres Pro Shardman extensions to the CREATE SEQUENCE command enable creation of global sequence number generators. This command creates an ordinary PostgreSQL sequence on all nodes in a cluster and records sequence parameters in the global sequence state dictionary.

After a global sequence is created, usual nextval function can be used to generate next sequence values that are guaranteed be unique across the entire cluster. Other standard sequence manipulation functions (e.g. setval) must not be used on global sequences as this may lead to unexpected results.

Although you cannot update a sequence directly, you can use a query like:

```
SELECT * FROM name;
```

to examine the parameters and current state of a sequence. In particular, the last_value field of the sequence shows the last value allocated by any session. (Of course, this value might be obsolete by the time it's printed, if other sessions are actively doing nextval calls.)

Parameters

TEMPORARY or TEMP

If specified, the sequence object is created only for this session, and is automatically dropped on session exit. Existing permanent sequences with the same name are not visible (in this session) while the temporary sequence exists, unless they are referenced with schema-qualified names.

UNLOGGED

If specified, the sequence is created as an unlogged sequence. Changes to unlogged sequences are not written to the write-ahead log. They are not crash-safe: an unlogged sequence is automatically

reset to its initial state after a crash or unclean shutdown. Unlogged sequences are also not replicated to standby servers.

Unlike unlogged tables, unlogged sequences do not offer a significant performance advantage. This option is mainly intended for sequences associated with unlogged tables via identity columns or serial columns. In those cases, it usually wouldn't make sense to have the sequence WAL-logged and replicated but not its associated table.

IF NOT EXISTS

Do not throw an error if a relation with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing relation is anything like the sequence that would have been created — it might not even be a sequence.

name

The name (optionally schema-qualified) of the sequence to be created.

data_type

The optional clause `AS data_type` specifies the data type of the sequence. Valid types are `smallint`, `integer`, and `bigint`. `bigint` is the default. The data type determines the default minimum and maximum values of the sequence.

increment

The optional clause `INCREMENT BY increment` specifies which value is added to the current sequence value to create a new value. A positive value will make an ascending sequence, a negative one a descending sequence. The default value is 1.

minvalue

NO MINVALUE

The optional clause `MINVALUE minvalue` determines the minimum value a sequence can generate. If this clause is not supplied or `NO MINVALUE` is specified, then defaults will be used. The default for an ascending sequence is 1. The default for a descending sequence is the minimum value of the data type.

maxvalue

NO MAXVALUE

The optional clause `MAXVALUE maxvalue` determines the maximum value for the sequence. If this clause is not supplied or `NO MAXVALUE` is specified, then default values will be used. The default for an ascending sequence is the maximum value of the data type. The default for a descending sequence is -1.

start

The optional clause `START WITH start` allows the sequence to begin anywhere. The default starting value is *minvalue* for ascending sequences and *maxvalue* for descending ones.

cache

The optional clause `CACHE cache` specifies how many sequence numbers are to be preallocated and stored in memory for faster access. The minimum value is 1 (only one value can be generated at a time, i.e., no cache), and this is also the default.

CYCLE

NO CYCLE

The `CYCLE` option allows the sequence to wrap around when the *maxvalue* or *minvalue* has been reached by an ascending or descending sequence respectively. If the limit is reached, the next number generated will be the *minvalue* or *maxvalue*, respectively.

If `NO CYCLE` is specified, any calls to `nextval` after the sequence has reached its maximum value will return an error. If neither `CYCLE` or `NO CYCLE` are specified, `NO CYCLE` is the default.

`OWNED BY table_name.column_name`
`OWNED BY NONE`

The `OWNED BY` option causes the sequence to be associated with a specific table column, such that if that column (or its whole table) is dropped, the sequence will be automatically dropped as well. The specified table must have the same owner and be in the same schema as the sequence. `OWNED BY NONE`, the default, specifies that there is no such association.

`global`

If specified, the sequence object is created as a Postgres Pro Shardman-managed global sequence.

`block_size`

The number of elements allocated for a local sequence. The default value is 65536.

Notes

Use `DROP SEQUENCE` to remove a sequence.

Sequences are based on `bigint` arithmetic, so the range cannot exceed the range of an eight-byte integer (-9223372036854775808 to 9223372036854775807).

Because `nextval` and `setval` calls are never rolled back, sequence objects cannot be used if “gapless” assignment of sequence numbers is needed. It is possible to build gapless assignment by using exclusive locking of a table containing a counter; but this solution is much more expensive than sequence objects, especially if many transactions need sequence numbers concurrently.

Unexpected results might be obtained if a `cache` setting greater than one is used for a sequence object that will be used concurrently by multiple sessions. Each session will allocate and cache successive sequence values during one access to the sequence object and increase the sequence object's `last_value` accordingly. Then, the next `cache-1` uses of `nextval` within that session simply return the preallocated values without touching the sequence object. So, any numbers allocated but not used within a session will be lost when that session ends, resulting in “holes” in the sequence.

Furthermore, although multiple sessions are guaranteed to allocate distinct sequence values, the values might be generated out of sequence when all the sessions are considered. For example, with a `cache` setting of 10, session A might reserve values 1..10 and return `nextval=1`, then session B might reserve values 11..20 and return `nextval=11` before session A has generated `nextval=2`. Thus, with a `cache` setting of one it is safe to assume that `nextval` values are generated sequentially; with a `cache` setting greater than one you should only assume that the `nextval` values are all distinct, not that they are generated purely sequentially. Also, `last_value` will reflect the latest value reserved by any session, whether or not it has yet been returned by `nextval`.

Another consideration is that a `setval` executed on such a sequence will not be noticed by other sessions until they have used up any preallocated values they have cached.

Examples

Create a global sequence called `serial`.

```
CREATE SEQUENCE serial MINVALUE 100 WITH (global);
```

Select the next number from this sequence:

```
SELECT nextval('serial');
 nextval
-----
      100
```

(1 row)

Create an ascending sequence called `serial`, starting at 101:

```
CREATE SEQUENCE serial START 101;
```

Select the next number from this sequence:

```
SELECT nextval('serial');
```

```
nextval
-----
      101
```

Select the next number from this sequence:

```
SELECT nextval('serial');
```

```
nextval
-----
      102
```

Use this sequence in an `INSERT` command:

```
INSERT INTO distributors VALUES (nextval('serial'), 'nothing');
```

Update the sequence value after a `COPY FROM`:

```
BEGIN;
COPY distributors FROM 'input_file';
SELECT setval('serial', max(id)) FROM distributors;
END;
```

Compatibility

`CREATE SEQUENCE` conforms to the SQL standard, with the following exceptions:

- Obtaining the next value is done using the `nextval()` function instead of the standard's `NEXT VALUE FOR expression`.
- The `OWNED BY` clause is a Postgres Pro extension.

See Also

[ALTER SEQUENCE](#), [DROP SEQUENCE](#)

CREATE SERVER

CREATE SERVER — define a new foreign server

Synopsis

```
CREATE SERVER [ IF NOT EXISTS ] server_name [ TYPE 'server_type' ] [ VERSION 'server_version' ]  
    FOREIGN DATA WRAPPER fdw_name  
    [ OPTIONS ( option 'value' [, ... ] ) ]
```

Description

CREATE SERVER defines a new foreign server. The user who defines the server becomes its owner.

A foreign server typically encapsulates connection information that a foreign-data wrapper uses to access an external data resource. Additional user-specific connection information may be specified by means of user mappings.

The server name must be unique within the database.

Creating a server requires USAGE privilege on the foreign-data wrapper being used.

Parameters

IF NOT EXISTS

Do not throw an error if a server with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing server is anything like the one that would have been created.

server_name

The name of the foreign server to be created.

server_type

Optional server type, potentially useful to foreign-data wrappers.

server_version

Optional server version, potentially useful to foreign-data wrappers.

fdw_name

The name of the foreign-data wrapper that manages the server.

OPTIONS (*option* '*value*' [, ...])

This clause specifies the options for the server. The options typically define the connection details of the server, but the actual names and values are dependent on the server's foreign-data wrapper.

Notes

When using the [dblink](#) module, a foreign server's name can be used as an argument of the [dblink_connect](#) function to indicate the connection parameters. It is necessary to have the USAGE privilege on the foreign server to be able to use it in this way.

If the foreign server supports sort pushdown, it is necessary for it to have the same sort ordering as the local server.

Examples

Create a server `myserver` that uses the foreign-data wrapper `postgres_fdw`:

```
CREATE SERVER myserver FOREIGN DATA WRAPPER postgres_fdw OPTIONS (host 'foo', dbname 'foodb', port '5432');
```

See [postgres_fdw](#) for more details.

Compatibility

`CREATE SERVER` conforms to ISO/IEC 9075-9 (SQL/MED).

See Also

[ALTER SERVER](#), [DROP SERVER](#)

CREATE STATISTICS

CREATE STATISTICS — define extended statistics

Synopsis

```
CREATE STATISTICS [ [ IF NOT EXISTS ] statistics_name ]
    ON ( expression )
    FROM table_name

CREATE STATISTICS [ [ IF NOT EXISTS ] statistics_name ]
    [ ( statistics_kind [, ...] ) ]
    ON { column_name | ( expression ) }, { column_name | ( expression ) } [, ...]
    FROM table_name
```

Description

CREATE STATISTICS will create a new extended statistics object tracking data about the specified table, foreign table or materialized view. The statistics object will be created in the current database and will be owned by the user issuing the command.

The CREATE STATISTICS command has two basic forms. The first form allows univariate statistics for a single expression to be collected, providing benefits similar to an expression index without the overhead of index maintenance. This form does not allow the statistics kind to be specified, since the various statistics kinds refer only to multivariate statistics. The second form of the command allows multivariate statistics on multiple columns and/or expressions to be collected, optionally specifying which statistics kinds to include. This form will also automatically cause univariate statistics to be collected on any expressions included in the list.

If a schema name is given (for example, CREATE STATISTICS myschema.mystat ...) then the statistics object is created in the specified schema. Otherwise it is created in the current schema. If given, the name of the statistics object must be distinct from the name of any other statistics object in the same schema.

Parameters

IF NOT EXISTS

Do not throw an error if a statistics object with the same name already exists. A notice is issued in this case. Note that only the name of the statistics object is considered here, not the details of its definition. Statistics name is required when IF NOT EXISTS is specified.

statistics_name

The name (optionally schema-qualified) of the statistics object to be created. If the name is omitted, Postgres Pro chooses a suitable name based on the parent table's name and the defined column name(s) and/or expression(s).

statistics_kind

A multivariate statistics kind to be computed in this statistics object. Currently supported kinds are `ndistinct`, which enables n-distinct statistics, `dependencies`, which enables functional dependency statistics, and `mcv` which enables most-common values lists. If this clause is omitted, all supported statistics kinds are included in the statistics object. Univariate expression statistics are built automatically if the statistics definition includes any complex expressions rather than just simple column references. For more information, see [Section 14.2.2](#) and [Section 68.2](#).

column_name

The name of a table column to be covered by the computed statistics. This is only allowed when building multivariate statistics. At least two column names or expressions must be specified, and their order is not significant.

expression

An expression to be covered by the computed statistics. This may be used to build univariate statistics on a single expression, or as part of a list of multiple column names and/or expressions to build multivariate statistics. In the latter case, separate univariate statistics are built automatically for each expression in the list.

table_name

The name (optionally schema-qualified) of the table containing the column(s) the statistics are computed on; see [ANALYZE](#) for an explanation of the handling of inheritance and partitions.

Notes

You must be the owner of a table to create a statistics object reading it. Once created, however, the ownership of the statistics object is independent of the underlying table(s).

Expression statistics are per-expression and are similar to creating an index on the expression, except that they avoid the overhead of index maintenance. Expression statistics are built automatically for each expression in the statistics object definition.

Extended statistics are not currently used by the planner for selectivity estimations made for table joins. This limitation will likely be removed in a future version of Postgres Pro.

Examples

Create table `t1` with two functionally dependent columns, i.e., knowledge of a value in the first column is sufficient for determining the value in the other column. Then functional dependency statistics are built on those columns:

```
CREATE TABLE t1 (  
    a    int,  
    b    int  
);  
  
INSERT INTO t1 SELECT i/100, i/500  
    FROM generate_series(1,1000000) s(i);  
  
ANALYZE t1;  
  
-- the number of matching rows will be drastically underestimated:  
EXPLAIN ANALYZE SELECT * FROM t1 WHERE (a = 1) AND (b = 0);  
  
CREATE STATISTICS s1 (dependencies) ON a, b FROM t1;  
  
ANALYZE t1;  
  
-- now the row count estimate is more accurate:  
EXPLAIN ANALYZE SELECT * FROM t1 WHERE (a = 1) AND (b = 0);
```

Without functional-dependency statistics, the planner would assume that the two `WHERE` conditions are independent, and would multiply their selectivities together to arrive at a much-too-small row count estimate. With such statistics, the planner recognizes that the `WHERE` conditions are redundant and does not underestimate the row count.

Create table `t2` with two perfectly correlated columns (containing identical data), and an MCV list on those columns:

```
CREATE TABLE t2 (  
    a    int,  
    b    int  
);
```

```
INSERT INTO t2 SELECT mod(i,100), mod(i,100)
      FROM generate_series(1,1000000) s(i);
```

```
CREATE STATISTICS s2 (mcv) ON a, b FROM t2;
```

```
ANALYZE t2;
```

```
-- valid combination (found in MCV)
EXPLAIN ANALYZE SELECT * FROM t2 WHERE (a = 1) AND (b = 1);
```

```
-- invalid combination (not found in MCV)
EXPLAIN ANALYZE SELECT * FROM t2 WHERE (a = 1) AND (b = 2);
```

The MCV list gives the planner more detailed information about the specific values that commonly appear in the table, as well as an upper bound on the selectivities of combinations of values that do not appear in the table, allowing it to generate better estimates in both cases.

Create table `t3` with a single timestamp column, and run queries using expressions on that column. Without extended statistics, the planner has no information about the data distribution for the expressions, and uses default estimates. The planner also does not realize that the value of the date truncated to the month is fully determined by the value of the date truncated to the day. Then expression and `ndistinct` statistics are built on those two expressions:

```
CREATE TABLE t3 (
    a    timestamp
);
```

```
INSERT INTO t3 SELECT i FROM generate_series('2020-01-01'::timestamp,
      '2020-12-31'::timestamp,
      '1 minute'::interval) s(i);
```

```
ANALYZE t3;
```

```
-- the number of matching rows will be drastically underestimated:
EXPLAIN ANALYZE SELECT * FROM t3
    WHERE date_trunc('month', a) = '2020-01-01'::timestamp;
```

```
EXPLAIN ANALYZE SELECT * FROM t3
    WHERE date_trunc('day', a) BETWEEN '2020-01-01'::timestamp
      AND '2020-06-30'::timestamp;
```

```
EXPLAIN ANALYZE SELECT date_trunc('month', a), date_trunc('day', a)
    FROM t3 GROUP BY 1, 2;
```

```
-- build ndistinct statistics on the pair of expressions (per-expression
-- statistics are built automatically)
CREATE STATISTICS s3 (ndistinct) ON date_trunc('month', a), date_trunc('day', a) FROM
    t3;
```

```
ANALYZE t3;
```

```
-- now the row count estimates are more accurate:
EXPLAIN ANALYZE SELECT * FROM t3
    WHERE date_trunc('month', a) = '2020-01-01'::timestamp;
```

```
EXPLAIN ANALYZE SELECT * FROM t3
    WHERE date_trunc('day', a) BETWEEN '2020-01-01'::timestamp
      AND '2020-06-30'::timestamp;
```

```
EXPLAIN ANALYZE SELECT date_trunc('month', a), date_trunc('day', a)
FROM t3 GROUP BY 1, 2;
```

Without expression and ndistinct statistics, the planner has no information about the number of distinct values for the expressions, and has to rely on default estimates. The equality and range conditions are assumed to have 0.5% selectivity, and the number of distinct values in the expression is assumed to be the same as for the column (i.e. unique). This results in a significant underestimate of the row count in the first two queries. Moreover, the planner has no information about the relationship between the expressions, so it assumes the two `WHERE` and `GROUP BY` conditions are independent, and multiplies their selectivities together to arrive at a severe overestimate of the group count in the aggregate query. This is further exacerbated by the lack of accurate statistics for the expressions, forcing the planner to use a default ndistinct estimate for the expression derived from ndistinct for the column. With such statistics, the planner recognizes that the conditions are correlated, and arrives at much more accurate estimates.

Compatibility

There is no `CREATE STATISTICS` command in the SQL standard.

See Also

[ALTER STATISTICS](#), [DROP STATISTICS](#)

CREATE SUBSCRIPTION

CREATE SUBSCRIPTION — define a new subscription

Synopsis

```
CREATE SUBSCRIPTION subscription_name
    CONNECTION 'conninfo'
    PUBLICATION publication_name [, ...]
    [ WITH ( subscription_parameter [= value] [, ... ] ) ]
```

Description

CREATE SUBSCRIPTION adds a new logical-replication subscription. The user that creates a subscription becomes the owner of the subscription. The subscription name must be distinct from the name of any existing subscription in the current database.

A subscription represents a replication connection to the publisher. Hence, in addition to adding definitions in the local catalogs, this command normally creates a replication slot on the publisher.

A logical replication worker will be started to replicate data for the new subscription at the commit of the transaction where this command is run, unless the subscription is initially disabled.

To be able to create a subscription, you must have the privileges of the `pg_create_subscription` role, as well as CREATE privileges on the current database.

Additional information about subscriptions and logical replication as a whole is available at [Section 29.2](#) and [Chapter 29](#).

Parameters

subscription_name

The name of the new subscription.

CONNECTION '*conninfo*'

The libpq connection string defining how to connect to the publisher database. For details see [Section 35.1.1](#).

PUBLICATION *publication_name* [, ...]

Names of the publications on the publisher to subscribe to.

WITH (*subscription_parameter* [= *value*] [, ...])

This clause specifies optional parameters for a subscription.

The following parameters control what happens during subscription creation:

connect (boolean)

Specifies whether the CREATE SUBSCRIPTION command should connect to the publisher at all. The default is true. Setting this to false will force the values of create_slot, enabled and copy_data to false. (You cannot combine setting connect to false with setting create_slot, enabled, or copy_data to true.)

Since no connection is made when this option is false, no tables are subscribed. To initiate replication, you must manually create the replication slot, enable the failover if required, enable the subscription, and refresh the subscription. See [Section 29.2.3](#) for examples.

`create_slot` (boolean)

Specifies whether the command should create the replication slot on the publisher. The default is `true`.

If set to `false`, you are responsible for creating the publisher's slot in some other way. See [Section 29.2.3](#) for examples.

`enabled` (boolean)

Specifies whether the subscription should be actively replicating or whether it should just be set up but not started yet. The default is `true`.

`slot_name` (string)

Name of the publisher's replication slot to use. The default is to use the name of the subscription for the slot name.

Setting `slot_name` to `NONE` means there will be no replication slot associated with the subscription. Such subscriptions must also have both `enabled` and `create_slot` set to `false`. Use this when you will be creating the replication slot later manually. See [Section 29.2.3](#) for examples.

When setting `slot_name` to a valid name and `create_slot` to `false`, the `failover` property value of the named slot may differ from the counterpart `failover` parameter specified in the subscription. Always ensure the slot property `failover` matches the counterpart parameter of the subscription and vice versa. Otherwise, the slot on the publisher may behave differently from what these subscription options say: for example, the slot on the publisher could either be synced to the standbys even when the subscription's `failover` option is disabled or could be disabled for sync even when the subscription's `failover` option is enabled.

The following parameters control the subscription's replication behavior after it has been created:

`binary` (boolean)

Specifies whether the subscription will request the publisher to send the data in binary format (as opposed to text). The default is `false`. Any initial table synchronization copy (see `copy_data`) also uses the same format. Binary format can be faster than the text format, but it is less portable across machine architectures and Postgres Pro versions. Binary format is very data type specific; for example, it will not allow copying from a `smallint` column to an `integer` column, even though that would work fine in text format. Even when this option is enabled, only data types having binary send and receive functions will be transferred in binary. Note that the initial synchronization requires all data types to have binary send and receive functions, otherwise the synchronization will fail (see [CREATE TYPE](#) for more about send/receive functions).

When doing cross-version replication, it could be that the publisher has a binary send function for some data type, but the subscriber lacks a binary receive function for that type. In such a case, data transfer will fail, and the `binary` option cannot be used.

If the publisher is a Postgres Pro version before 16, then any initial table synchronization will use text format even if `binary = true`.

`copy_data` (boolean)

Specifies whether to copy pre-existing data in the publications that are being subscribed to when the replication starts. The default is `true`.

If the publications contain `WHERE` clauses, it will affect what data is copied. Refer to the [Notes](#) for details.

See [Notes](#) for details of how `copy_data = true` can interact with the `origin` parameter.

streaming (enum)

Specifies whether to enable streaming of in-progress transactions for this subscription. The default value is `off`, meaning all transactions are fully decoded on the publisher and only then sent to the subscriber as a whole.

If set to `on`, the incoming changes are written to temporary files and then applied only after the transaction is committed on the publisher and received by the subscriber.

If set to `parallel`, incoming changes are directly applied via one of the parallel apply workers, if available. If no parallel apply worker is free to handle streaming transactions then the changes are written to temporary files and applied after the transaction is committed. Note that if an error happens in a parallel apply worker, the finish LSN of the remote transaction might not be reported in the server log.

synchronous_commit (enum)

The value of this parameter overrides the [synchronous_commit](#) setting within this subscription's apply worker processes. The default value is `off`.

It is safe to use `off` for logical replication: If the subscriber loses transactions because of missing synchronization, the data will be sent again from the publisher.

A different setting might be appropriate when doing synchronous logical replication. The logical replication workers report the positions of writes and flushes to the publisher, and when using synchronous replication, the publisher will wait for the actual flush. This means that setting `synchronous_commit` for the subscriber to `off` when the subscription is used for synchronous replication might increase the latency for `COMMIT` on the publisher. In this scenario, it can be advantageous to set `synchronous_commit` to `local` or higher.

two_phase (boolean)

Specifies whether two-phase commit is enabled for this subscription. The default is `false`.

When two-phase commit is enabled, prepared transactions are sent to the subscriber at the time of `PREPARE TRANSACTION`, and are processed as two-phase transactions on the subscriber too. Otherwise, prepared transactions are sent to the subscriber only when committed, and are then processed immediately by the subscriber.

The implementation of two-phase commit requires that replication has successfully finished the initial table synchronization phase. So even when `two_phase` is enabled for a subscription, the internal two-phase state remains temporarily “pending” until the initialization phase completes. See column `subtwophasestate` of [pg_subscription](#) to know the actual two-phase state.

disable_on_error (boolean)

Specifies whether the subscription should be automatically disabled if any errors are detected by subscription workers during data replication from the publisher. The default is `false`.

password_required (boolean)

If set to `true`, connections to the publisher made as a result of this subscription must use password authentication and the password must be specified as a part of the connection string. This setting is ignored when the subscription is owned by a superuser. The default is `true`. Only superusers can set this value to `false`.

run_as_owner (boolean)

If `true`, all replication actions are performed as the subscription owner. If `false`, replication workers will perform actions on each table as the owner of that table. The latter configuration is generally much more secure; for details, see [Section 29.10](#). The default is `false`.

`origin (string)`

Specifies whether the subscription will request the publisher to only send changes that don't have an origin or send changes regardless of origin. Setting `origin` to `none` means that the subscription will request the publisher to only send changes that don't have an origin. Setting `origin` to `any` means that the publisher sends changes regardless of their origin. The default is `any`.

See [Notes](#) for details of how `copy_data = true` can interact with the `origin` parameter.

`failover (boolean)`

Specifies whether the replication slots associated with the subscription are enabled to be synced to the standbys so that logical replication can be resumed from the new primary after failover. The default is `false`.

When specifying a parameter of type `boolean`, the `= value` part can be omitted, which is equivalent to specifying `TRUE`.

Notes

See [Section 29.10](#) for details on how to configure access control between the subscription and the publication instance.

When creating a replication slot (the default behavior), `CREATE SUBSCRIPTION` cannot be executed inside a transaction block.

Creating a subscription that connects to the same database cluster (for example, to replicate between databases in the same cluster or to replicate within the same database) will only succeed if the replication slot is not created as part of the same command. Otherwise, the `CREATE SUBSCRIPTION` call will hang. To make this work, create the replication slot separately (using the function `pg_create_logical_replication_slot` with the plugin name `pgoutput`) and create the subscription using the parameter `create_slot = false`. See [Section 29.2.3](#) for examples. This is an implementation restriction that might be lifted in a future release.

If any table in the publication has a `WHERE` clause, rows for which the *expression* evaluates to false or null will not be published. If the subscription has several publications in which the same table has been published with different `WHERE` clauses, a row will be published if any of the expressions (referring to that publish operation) are satisfied. In the case of different `WHERE` clauses, if one of the publications has no `WHERE` clause (referring to that publish operation) or the publication is declared as [FOR ALL TABLES](#) or [FOR TABLES IN SCHEMA](#), rows are always published regardless of the definition of the other expressions. If the subscriber is a Postgres Pro version before 15, then any row filtering is ignored during the initial data synchronization phase. For this case, the user might want to consider deleting any initially copied data that would be incompatible with subsequent filtering. Because initial data synchronization does not take into account the publication `publish` parameter when copying existing table data, some rows may be copied that would not be replicated using DML. See [Section 29.2.2](#) for examples.

Subscriptions having several publications in which the same table has been published with different column lists are not supported.

We allow non-existent publications to be specified so that users can add those later. This means `pg_subscription` can have non-existent publications.

When using a subscription parameter combination of `copy_data = true` and `origin = NONE`, the initial sync table data is copied directly from the publisher, meaning that knowledge of the true origin of that data is not possible. If the publisher also has subscriptions then the copied table data might have originated from further upstream. This scenario is detected and a `WARNING` is logged to the user, but the warning is only an indication of a potential problem; it is the user's responsibility to make the necessary checks to ensure the copied data origins are really as wanted or not.

To find which tables might potentially include non-local origins (due to other subscriptions created on the publisher) try this SQL query:

```
# substitute <pub-names> below with your publication name(s) to be queried
SELECT DISTINCT PT.schemaname, PT.tablename
FROM pg_publication_tables PT
     JOIN pg_class C ON (C.relname = PT.tablename)
     JOIN pg_namespace N ON (N.nspname = PT.schemaname),
     pg_subscription_rel PS
WHERE C.relnamespace = N.oid AND
      (PS.srrelid = C.oid OR
       C.oid IN (SELECT relid FROM pg_partition_ancestors(PS.srrelid) UNION
                SELECT relid FROM pg_partition_tree(PS.srrelid))) AND
      PT.pubname IN (<pub-names>);
```

Examples

Create a subscription to a remote server that replicates tables in the publications `mypublication` and `insert_only` and starts replicating immediately on commit:

```
CREATE SUBSCRIPTION mysub
    CONNECTION 'host=192.168.1.50 port=5432 user=foo dbname=foodb'
    PUBLICATION mypublication, insert_only;
```

Create a subscription to a remote server that replicates tables in the `insert_only` publication and does not start replicating until enabled at a later time.

```
CREATE SUBSCRIPTION mysub
    CONNECTION 'host=192.168.1.50 port=5432 user=foo dbname=foodb'
    PUBLICATION insert_only
    WITH (enabled = false);
```

Compatibility

`CREATE SUBSCRIPTION` is a Postgres Pro extension.

See Also

[ALTER SUBSCRIPTION](#), [DROP SUBSCRIPTION](#), [CREATE PUBLICATION](#), [ALTER PUBLICATION](#)

CREATE TABLE

CREATE TABLE — define a new table

Synopsis

For sharded tables:

```
CREATE TABLE [ IF NOT EXISTS ] sharded_table_name ( [
    { column_name data_type [ COLLATE collation ] [ column_constraint [ ... ] ]
    | sharded_table_constraint
    | LIKE source_table [ like_option ... ] }
    [, ... ]
] )
[ USING method ]
[ TABLESPACE tablespace_name ]
[ WITH ( { distributed_by = 'column_name [, ... ]'
    [, num_parts = number_of_partitions ]
    [, colocate_with = 'colocation_table_name' ]
    [, partition_by = 'column_name',
        partition_bounds = 'array_of_partition_bound_exprs' ] |
        global = true
    [, storage_parameter [= value] [, ... ] }
)]
```

where *column_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ NOT NULL |
  NULL |
  CHECK ( expression ) [ NO INHERIT ] |
  DEFAULT default_expr |
  UNIQUE index_parameters |
  PRIMARY KEY index_parameters |
  REFERENCES reftable [ ( refcolumn ) ] [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ]
    [ ON DELETE referential_action ] [ ON UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *sharded_table_constraint* is:

```
[ CONSTRAINT constraint_name ]
{ CHECK ( expression ) [ NO INHERIT ] |
  UNIQUE ( column_name [, ... ] ) index_parameters |
  PRIMARY KEY ( column_name [, ... ] ) index_parameters |
  EXCLUDE [ USING index_method ] ( exclude_element WITH operator
    [, ... ] ) index_parameters [ WHERE ( predicate ) ]
  FOREIGN KEY ( column_name [, ... ] ) REFERENCES reftable [ ( refcolumn [, ... ] ) ]
    [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ] [ ON DELETE referential_action ] [ ON
    UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

For other table types:

```
CREATE [ [ GLOBAL | LOCAL ] { TEMPORARY | TEMP } | UNLOGGED | CONSTANT ] TABLE [ IF NOT
EXISTS ] table_name ( [
    { column_name data_type [ STORAGE { PLAIN | EXTERNAL | EXTENDED | MAIN | DEFAULT } ]
    [ COMPRESSION compression_method ] [ COLLATE collation ] [ column_constraint [ ... ] ]
```

```

    | table_constraint
    | LIKE source_table [ like_option ... ] }
    [, ... ]
] )
[ INHERITS ( parent_table [, ... ] ) ]
[ PARTITION BY { RANGE | LIST | HASH } ( { column_name | ( expression ) }
  [ COLLATE collation ] [ opclass ] [, ... ] ) ]
[ USING method ]
[ WITH ( storage_parameter [= value] [, ... ] ) | WITHOUT OIDS ]
[ ON COMMIT { PRESERVE ROWS | DELETE ROWS | DROP } ]
[ TABLESPACE tablespace_name ]

CREATE [ [ GLOBAL | LOCAL ] { TEMPORARY | TEMP } | UNLOGGED | CONSTANT ] TABLE [ IF NOT
  EXISTS ] table_name
  OF type_name [ (
    { column_name [ WITH OPTIONS ] [ column_constraint [ ... ] ]
      | table_constraint }
    [, ... ]
  ) ]
[ PARTITION BY { RANGE | LIST | HASH } ( { column_name | ( expression ) }
  [ COLLATE collation ] [ opclass ] [, ... ] ) ]
[ USING method ]
[ WITH ( storage_parameter [= value] [, ... ] ) | WITHOUT OIDS ]
[ ON COMMIT { PRESERVE ROWS | DELETE ROWS | DROP } ]
[ TABLESPACE tablespace_name ]

CREATE [ [ GLOBAL | LOCAL ] { TEMPORARY | TEMP } | UNLOGGED ] TABLE [ IF NOT
  EXISTS ] table_name
  PARTITION OF parent_table [ (
    { column_name [ WITH OPTIONS ] [ column_constraint [ ... ] ]
      | table_constraint }
    [, ... ]
  ) ] { FOR VALUES partition_bound_spec | DEFAULT }
[ PARTITION BY { RANGE | LIST | HASH } ( { column_name | ( expression ) }
  [ COLLATE collation ] [ opclass ] [, ... ] ) ]
[ USING method ]
[ WITH ( storage_parameter [= value] [, ... ] ) | WITHOUT OIDS ]
[ ON COMMIT { PRESERVE ROWS | DELETE ROWS | DROP } ]
[ TABLESPACE tablespace_name ]

where column_constraint is:

[ CONSTRAINT constraint_name ]
{ NOT NULL |
  NULL |
  CHECK ( expression ) [ NO INHERIT ] |
  DEFAULT default_expr |
  GENERATED ALWAYS AS ( generation_expr ) STORED |
  GENERATED { ALWAYS | BY DEFAULT } AS IDENTITY [ ( sequence_options ) ] |
  UNIQUE [ NULLS [ NOT ] DISTINCT ] index_parameters |
  PRIMARY KEY index_parameters |
  REFERENCES reftable [ ( refcolumn ) ] [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ]
    [ ON DELETE referential_action ] [ ON UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]

and table_constraint is:

[ CONSTRAINT constraint_name ]

```

```
{ CHECK ( expression ) [ NO INHERIT ] |
  UNIQUE [ NULLS [ NOT ] DISTINCT ] ( column_name [, ... ] ) index_parameters |
  PRIMARY KEY ( column_name [, ... ] ) index_parameters |
  EXCLUDE [ USING index_method ] ( exclude_element WITH operator
[, ... ] ) index_parameters [ WHERE ( predicate ) ] |
  FOREIGN KEY ( column_name [, ... ] ) REFERENCES reftable [ ( refcolumn [, ... ] ) ]
    [ MATCH FULL | MATCH PARTIAL | MATCH SIMPLE ] [ ON DELETE referential_action ] [ ON
UPDATE referential_action ] }
[ DEFERRABLE | NOT DEFERRABLE ] [ INITIALLY DEFERRED | INITIALLY IMMEDIATE ]
```

and *like_option* is:

```
{ INCLUDING | EXCLUDING } { COMMENTS | COMPRESSION | CONSTRAINTS | DEFAULTS | GENERATED
| IDENTITY | INDEXES | STATISTICS | STORAGE | ALL }
```

and *partition_bound_spec* is:

```
IN ( partition_bound_expr [, ...] ) |
FROM ( { partition_bound_expr | MINVALUE | MAXVALUE } [, ...] )
  TO ( { partition_bound_expr | MINVALUE | MAXVALUE } [, ...] ) |
WITH ( MODULUS numeric_literal, REMAINDER numeric_literal )
```

index_parameters in UNIQUE, PRIMARY KEY, and EXCLUDE constraints are:

```
[ INCLUDE ( column_name [, ... ] ) ]
[ WITH ( storage_parameter [= value] [, ... ] ) ]
[ USING INDEX TABLESPACE tablespace_name ]
```

exclude_element in an EXCLUDE constraint is:

```
{ column_name | ( expression ) } [ COLLATE collation ] [ opclass [ ( opclass_parameter
= value [, ... ] ) ] ] [ ASC | DESC ] [ NULLS { FIRST | LAST } ]
```

referential_action in a FOREIGN KEY/REFERENCES constraint is:

```
{ NO ACTION | RESTRICT | CASCADE | SET NULL [ ( column_name [, ... ] ) ] | SET DEFAULT
[ ( column_name [, ... ] ) ] }
```

Description

CREATE TABLE will create a new, initially empty table in the current database. The table will be owned by the user issuing the command.

If a schema name is given (for example, CREATE TABLE myschema.mytable ...) then the table is created in the specified schema. Otherwise it is created in the current schema. Temporary tables exist in a special schema, so a schema name cannot be given when creating a temporary table. The name of the table must be distinct from the name of any other relation (table, sequence, index, view, materialized view, or foreign table) in the same schema.

CREATE TABLE also automatically creates a data type that represents the composite type corresponding to one row of the table. Therefore, tables cannot have the same name as any existing data type in the same schema.

The optional constraint clauses specify constraints (tests) that new or updated rows must satisfy for an insert or update operation to succeed. A constraint is an SQL object that helps define the set of valid values in the table in various ways.

There are two ways to define constraints: table constraints and column constraints. A column constraint is defined as part of a column definition. A table constraint definition is not tied to a particular column,

and it can encompass more than one column. Every column constraint can also be written as a table constraint; a column constraint is only a notational convenience for use when the constraint only affects one column.

To be able to create a table, you must have `USAGE` privilege on all column types or the type in the `OF` clause, respectively.

Parameters

TEMPORARY or TEMP

If specified, the table is created as a temporary table. Temporary tables are automatically dropped at the end of a session, or optionally at the end of the current transaction (see `ON COMMIT` below). The default `search_path` includes the temporary schema first and so identically named existing permanent tables are not chosen for new plans while the temporary table exists, unless they are referenced with schema-qualified names. Any indexes created on a temporary table are automatically temporary as well.

Unlike for permanent tables, we extend the physical file for a temporary table only when the table does not fit into `temp_buffers` cache anymore. Thus, the total and on-disk size of a temporary table can differ. Blocks are not allocated in advance, so we may run out of free space on disk when trying to evict a temporary buffer from cache. If the available disk space is smaller than `temp_buffers` size, an error occurs. This case is quite unusual, though.

The [autovacuum daemon](#) cannot access and therefore cannot vacuum or analyze temporary tables. For this reason, appropriate vacuum and analyze operations should be performed via session SQL commands. For example, if a temporary table is going to be used in complex queries, it is wise to run `ANALYZE` on the temporary table after it is populated.

Superusers can only use the `DROP TABLE` command with temporary relations of other sessions. Dropping such relations is not allowed when the `skip_temp_rel_lock` configuration parameter is enabled.

Optionally, `GLOBAL` or `LOCAL` can be written before `TEMPORARY` or `TEMP`. This presently makes no difference in Postgres Pro and is deprecated; see [Compatibility](#) below.

UNLOGGED

If specified, the table is created as an unlogged table. Data written to unlogged tables is not written to the write-ahead log (see [Chapter 28](#)), which makes them considerably faster than ordinary tables. However, they are not crash-safe: an unlogged table is automatically truncated after a crash or unclean shutdown. The contents of an unlogged table are also not replicated to standby servers. Any indexes created on an unlogged table are automatically unlogged as well.

If this is specified, any sequences created together with the unlogged table (for identity or serial columns) are also created as unlogged.

CONSTANT

If specified, the table is created as read-only. No data can be modified or added to constant tables, and they are not processed by [autovacuum](#). Constant tables cannot be changed to read-write mode, so there is no much sense to create them as constant; use [CREATE TABLE AS](#) or [ALTER TABLE](#) instead.

IF NOT EXISTS

Do not throw an error if a relation with the same name already exists. A notice is issued in this case. Note that there is no guarantee that the existing relation is anything like the one that would have been created.

table_name

The name (optionally schema-qualified) of the table to be created.

OF *type_name*

Creates a *typed table*, which takes its structure from the specified composite type (name optionally schema-qualified). A typed table is tied to its type; for example the table will be dropped if the type is dropped (with `DROP TYPE ... CASCADE`).

When a typed table is created, then the data types of the columns are determined by the underlying composite type and are not specified by the `CREATE TABLE` command. But the `CREATE TABLE` command can add defaults and constraints to the table and can specify storage parameters.

column_name

The name of a column to be created in the new table.

data_type

The data type of the column. This can include array specifiers. For more information on the data types supported by Postgres Pro, refer to [Chapter 8](#).

`COLLATE collation`

The `COLLATE` clause assigns a collation to the column (which must be of a collatable data type). If not specified, the column data type's default collation is used.

`STORAGE { PLAIN | EXTERNAL | EXTENDED | MAIN | DEFAULT }`

This form sets the storage mode for the column. This controls whether this column is held inline or in a secondary TOAST table, and whether the data should be compressed or not. `PLAIN` must be used for fixed-length values such as `integer` and is inline, uncompressed. `MAIN` is for inline, compressible data. `EXTERNAL` is for external, uncompressed data, and `EXTENDED` is for external, compressed data. Writing `DEFAULT` sets the storage mode to the default mode for the column's data type. `EXTENDED` is the default for most data types that support non-`PLAIN` storage. Use of `EXTERNAL` will make substring operations on very large `text` and `bytea` values run faster, at the penalty of increased storage space. See [Section 66.2](#) for more information.

`COMPRESSION compression_method`

The `COMPRESSION` clause sets the compression method for the column. Compression is supported only for variable-width data types, and is used only when the column's storage mode is `main` or `extended`. (See [ALTER TABLE](#) for information on column storage modes.) Setting this property for a partitioned table has no direct effect, because such tables have no storage of their own, but the configured value will be inherited by newly-created partitions. The supported compression methods are `pglz` and `lz4`. (`lz4` is available only if `--with-lz4` was used when building Postgres Pro.) In addition, *compression_method* can be `default` to explicitly specify the default behavior, which is to consult the [default_toast_compression](#) setting at the time of data insertion to determine the method to use.

`INHERITS (parent_table [, ...])`

The optional `INHERITS` clause specifies a list of tables from which the new table automatically inherits all columns. Parent tables can be plain tables or foreign tables.

Use of `INHERITS` creates a persistent relationship between the new child table and its parent table(s). Schema modifications to the parent(s) normally propagate to children as well, and by default the data of the child table is included in scans of the parent(s).

If the same column name exists in more than one parent table, an error is reported unless the data types of the columns match in each of the parent tables. If there is no conflict, then the duplicate columns are merged to form a single column in the new table. If the column name list of the new table contains a column name that is also inherited, the data type must likewise match the inherited column(s), and the column definitions are merged into one. If the new table explicitly specifies a default value for the column, this default overrides any defaults from inherited declarations of the column. Otherwise, any parents that specify default values for the column must all specify the same default, or an error will be reported.

CHECK constraints are merged in essentially the same way as columns: if multiple parent tables and/or the new table definition contain identically-named CHECK constraints, these constraints must all have the same check expression, or an error will be reported. Constraints having the same name and expression will be merged into one copy. A constraint marked NO INHERIT in a parent will not be considered. Notice that an unnamed CHECK constraint in the new table will never be merged, since a unique name will always be chosen for it.

Column STORAGE settings are also copied from parent tables.

If a column in the parent table is an identity column, that property is not inherited. A column in the child table can be declared identity column if desired.

```
PARTITION BY { RANGE | LIST | HASH } ( { column_name | ( expression ) } [ opclass ] [, ...] )
```

The optional PARTITION BY clause specifies a strategy of partitioning the table. The table thus created is called a *partitioned* table. The parenthesized list of columns or expressions forms the *partition key* for the table. When using range or hash partitioning, the partition key can include multiple columns or expressions (up to 32, but this limit can be altered when building Postgres Pro), but for list partitioning, the partition key must consist of a single column or expression.

Range and list partitioning require a btree operator class, while hash partitioning requires a hash operator class. If no operator class is specified explicitly, the default operator class of the appropriate type will be used; if no default operator class exists, an error will be raised. When hash partitioning is used, the operator class used must implement support function 2 (see [Section 39.16.3](#) for details).

A partitioned table is divided into sub-tables (called partitions), which are created using separate CREATE TABLE commands. The partitioned table is itself empty. A data row inserted into the table is routed to a partition based on the value of columns or expressions in the partition key. If no existing partition matches the values in the new row, an error will be reported.

Partitioned tables do not support EXCLUDE constraints; however, you can define these constraints on individual partitions.

See [Section 5.12](#) for more discussion on table partitioning.

```
PARTITION OF parent_table { FOR VALUES partition_bound_spec | DEFAULT }
```

Creates the table as a *partition* of the specified parent table. The table can be created either as a partition for specific values using FOR VALUES or as a default partition using DEFAULT. Any indexes, constraints and user-defined row-level triggers that exist in the parent table are cloned on the new partition.

The *partition_bound_spec* must correspond to the partitioning method and partition key of the parent table, and must not overlap with any existing partition of that parent. The form with IN is used for list partitioning, the form with FROM and TO is used for range partitioning, and the form with WITH is used for hash partitioning.

partition_bound_expr is any variable-free expression (subqueries, window functions, aggregate functions, and set-returning functions are not allowed). Its data type must match the data type of the corresponding partition key column. The expression is evaluated once at table creation time, so it can even contain volatile expressions such as CURRENT_TIMESTAMP.

When creating a list partition, NULL can be specified to signify that the partition allows the partition key column to be null. However, there cannot be more than one such list partition for a given parent table. NULL cannot be specified for range partitions.

When creating a range partition, the lower bound specified with FROM is an inclusive bound, whereas the upper bound specified with TO is an exclusive bound. That is, the values specified in the FROM list are valid values of the corresponding partition key columns for this partition, whereas those in the TO list are not. Note that this statement must be understood according to the rules of row-wise comparison ([Section 9.25.5](#)). For example, given PARTITION BY RANGE (x,y), a partition bound FROM (1, 2) TO (3, 4) allows x=1 with any y>=2, x=2 with any non-null y, and x=3 with any y<4.

The special values `MINVALUE` and `MAXVALUE` may be used when creating a range partition to indicate that there is no lower or upper bound on the column's value. For example, a partition defined using `FROM (MINVALUE) TO (10)` allows any values less than 10, and a partition defined using `FROM (10) TO (MAXVALUE)` allows any values greater than or equal to 10.

When creating a range partition involving more than one column, it can also make sense to use `MAXVALUE` as part of the lower bound, and `MINVALUE` as part of the upper bound. For example, a partition defined using `FROM (0, MAXVALUE) TO (10, MAXVALUE)` allows any rows where the first partition key column is greater than 0 and less than or equal to 10. Similarly, a partition defined using `FROM ('a', MINVALUE) TO ('b', MINVALUE)` allows any rows where the first partition key column starts with "a".

Note that if `MINVALUE` or `MAXVALUE` is used for one column of a partitioning bound, the same value must be used for all subsequent columns. For example, `(10, MINVALUE, 0)` is not a valid bound; you should write `(10, MINVALUE, MINVALUE)`.

Also note that some element types, such as `timestamp`, have a notion of "infinity", which is just another value that can be stored. This is different from `MINVALUE` and `MAXVALUE`, which are not real values that can be stored, but rather they are ways of saying that the value is unbounded. `MAXVALUE` can be thought of as being greater than any other value, including "infinity" and `MINVALUE` as being less than any other value, including "minus infinity". Thus the range `FROM ('infinity') TO (MAXVALUE)` is not an empty range; it allows precisely one value to be stored — "infinity".

If `DEFAULT` is specified, the table will be created as the default partition of the parent table. This option is not available for hash-partitioned tables. A partition key value not fitting into any other partition of the given parent will be routed to the default partition.

When a table has an existing `DEFAULT` partition and a new partition is added to it, the default partition must be scanned to verify that it does not contain any rows which properly belong in the new partition. If the default partition contains a large number of rows, this may be slow. The scan will be skipped if the default partition is a foreign table or if it has a constraint which proves that it cannot contain rows which should be placed in the new partition.

When creating a hash partition, a modulus and remainder must be specified. The modulus must be a positive integer, and the remainder must be a non-negative integer less than the modulus. Typically, when initially setting up a hash-partitioned table, you should choose a modulus equal to the number of partitions and assign every table the same modulus and a different remainder (see examples, below). However, it is not required that every partition have the same modulus, only that every modulus which occurs among the partitions of a hash-partitioned table is a factor of the next larger modulus. This allows the number of partitions to be increased incrementally without needing to move all the data at once. For example, suppose you have a hash-partitioned table with 8 partitions, each of which has modulus 8, but find it necessary to increase the number of partitions to 16. You can detach one of the modulus-8 partitions, create two new modulus-16 partitions covering the same portion of the key space (one with a remainder equal to the remainder of the detached partition, and the other with a remainder equal to that value plus 8), and repopulate them with data. You can then repeat this -- perhaps at a later time -- for each modulus-8 partition until none remain. While this may still involve a large amount of data movement at each step, it is still better than having to create a whole new table and move all the data at once.

A partition must have the same column names and types as the partitioned table to which it belongs. Modifications to the column names or types of a partitioned table will automatically propagate to all partitions. `CHECK` constraints will be inherited automatically by every partition, but an individual partition may specify additional `CHECK` constraints; additional constraints with the same name and condition as in the parent will be merged with the parent constraint. Defaults may be specified separately for each partition. But note that a partition's default value is not applied when inserting a tuple through a partitioned table.

Rows inserted into a partitioned table will be automatically routed to the correct partition. If no suitable partition exists, an error will occur.

Operations such as `TRUNCATE` which normally affect a table and all of its inheritance children will cascade to all partitions, but may also be performed on an individual partition.

Note that creating a partition using `PARTITION OF` requires taking an `ACCESS EXCLUSIVE` lock on the parent partitioned table. Likewise, dropping a partition with `DROP TABLE` requires taking an `ACCESS EXCLUSIVE` lock on the parent table. It is possible to use `ALTER TABLE ATTACH/DETACH PARTITION` to perform these operations with a weaker lock, thus reducing interference with concurrent operations on the partitioned table.

`LIKE source_table [like_option ...]`

The `LIKE` clause specifies a table from which the new table automatically copies all column names, their data types, and their not-null constraints.

Unlike `INHERITS`, the new table and original table are completely decoupled after creation is complete. Changes to the original table will not be applied to the new table, and it is not possible to include data of the new table in scans of the original table.

Also unlike `INHERITS`, columns and constraints copied by `LIKE` are not merged with similarly named columns and constraints. If the same name is specified explicitly or in another `LIKE` clause, an error is signaled.

The optional `like_option` clauses specify which additional properties of the original table to copy. Specifying `INCLUDING` copies the property, specifying `EXCLUDING` omits the property. `EXCLUDING` is the default. If multiple specifications are made for the same kind of object, the last one is used. The available options are:

`INCLUDING COMMENTS`

Comments for the copied columns, constraints, and indexes will be copied. The default behavior is to exclude comments, resulting in the copied columns and constraints in the new table having no comments.

`INCLUDING COMPRESSION`

Compression method of the columns will be copied. The default behavior is to exclude compression methods, resulting in columns having the default compression method.

`INCLUDING CONSTRAINTS`

`CHECK` constraints will be copied. No distinction is made between column constraints and table constraints. Not-null constraints are always copied to the new table.

`INCLUDING DEFAULTS`

Default expressions for the copied column definitions will be copied. Otherwise, default expressions are not copied, resulting in the copied columns in the new table having null defaults. Note that copying defaults that call database-modification functions, such as `nextval`, may create a functional linkage between the original and new tables.

`INCLUDING GENERATED`

Any generation expressions of copied column definitions will be copied. By default, new columns will be regular base columns.

`INCLUDING IDENTITY`

Any identity specifications of copied column definitions will be copied. A new sequence is created for each identity column of the new table, separate from the sequences associated with the old table.

`INCLUDING INDEXES`

Indexes, `PRIMARY KEY`, `UNIQUE`, and `EXCLUDE` constraints on the original table will be created on the new table. Names for the new indexes and constraints are chosen according to the default

rules, regardless of how the originals were named. (This behavior avoids possible duplicate-name failures for the new indexes.)

INCLUDING STATISTICS

Extended statistics are copied to the new table.

INCLUDING STORAGE

`STORAGE` settings for the copied column definitions will be copied. The default behavior is to exclude `STORAGE` settings, resulting in the copied columns in the new table having type-specific default settings. For more on `STORAGE` settings, see [Section 66.2](#).

INCLUDING ALL

`INCLUDING ALL` is an abbreviated form selecting all the available individual options. (It could be useful to write individual `EXCLUDING` clauses after `INCLUDING ALL` to select all but some specific options.)

The `LIKE` clause can also be used to copy column definitions from views, foreign tables, or composite types. Inapplicable options (e.g., `INCLUDING INDEXES` from a view) are ignored.

CONSTRAINT *constraint_name*

An optional name for a column or table constraint. If the constraint is violated, the constraint name is present in error messages, so constraint names like `col must be positive` can be used to communicate helpful constraint information to client applications. (Double-quotes are needed to specify constraint names that contain spaces.) If a constraint name is not specified, the system generates a name.

NOT NULL

The column is not allowed to contain null values.

NULL

The column is allowed to contain null values. This is the default.

This clause is only provided for compatibility with non-standard SQL databases. Its use is discouraged in new applications.

CHECK (*expression*) [NO INHERIT]

The `CHECK` clause specifies an expression producing a Boolean result which new or updated rows must satisfy for an insert or update operation to succeed. Expressions evaluating to `TRUE` or `UNKNOWN` succeed. Should any row of an insert or update operation produce a `FALSE` result, an error exception is raised and the insert or update does not alter the database. A check constraint specified as a column constraint should reference that column's value only, while an expression appearing in a table constraint can reference multiple columns.

Currently, `CHECK` expressions cannot contain subqueries nor refer to variables other than columns of the current row (see [Section 5.5.1](#)). The system column `tableoid` may be referenced, but not any other system column.

A constraint marked with `NO INHERIT` will not propagate to child tables.

When a table has multiple `CHECK` constraints, they will be tested for each row in alphabetical order by name, after checking `NOT NULL` constraints. (PostgreSQL versions before 9.5 did not honor any particular firing order for `CHECK` constraints.)

DEFAULT *default_expr*

The `DEFAULT` clause assigns a default data value for the column whose column definition it appears within. The value is any variable-free expression (in particular, cross-references to other columns in

the current table are not allowed). Subqueries are not allowed either. The data type of the default expression must match the data type of the column.

The default expression will be used in any insert operation that does not specify a value for the column. If there is no default for a column, then the default is null.

`GENERATED ALWAYS AS (generation_expr) STORED`

This clause creates the column as a *generated column*. The column cannot be written to, and when read the result of the specified expression will be returned.

The keyword `STORED` is required to signify that the column will be computed on write and will be stored on disk.

The generation expression can refer to other columns in the table, but not other generated columns. Any functions and operators used must be immutable. References to other tables are not allowed.

`GENERATED { ALWAYS | BY DEFAULT } AS IDENTITY [(sequence_options)]`

This clause creates the column as an *identity column*. It will have an implicit sequence attached to it and in newly-inserted rows the column will automatically have values from the sequence assigned to it. Such a column is implicitly `NOT NULL`.

The clauses `ALWAYS` and `BY DEFAULT` determine how explicitly user-specified values are handled in `INSERT` and `UPDATE` commands.

In an `INSERT` command, if `ALWAYS` is selected, a user-specified value is only accepted if the `INSERT` statement specifies `OVERRIDING SYSTEM VALUE`. If `BY DEFAULT` is selected, then the user-specified value takes precedence. See [INSERT](#) for details. (In the `COPY` command, user-specified values are always used regardless of this setting.)

In an `UPDATE` command, if `ALWAYS` is selected, any update of the column to any value other than `DEFAULT` will be rejected. If `BY DEFAULT` is selected, the column can be updated normally. (There is no `OVERRIDING` clause for the `UPDATE` command.)

The optional *sequence_options* clause can be used to override the parameters of the sequence. The available options include those shown for [CREATE SEQUENCE](#), plus `SEQUENCE NAME name`, `LOGGED`, and `UNLOGGED`, which allow selection of the name and persistence level of the sequence. Without `SEQUENCE NAME`, the system chooses an unused name for the sequence. Without `LOGGED` or `UNLOGGED`, the sequence will have the same persistence level as the table.

`UNIQUE [NULLS [NOT] DISTINCT] (column constraint)`

`UNIQUE [NULLS [NOT] DISTINCT] (column_name [, ...]) [INCLUDE (column_name [, ...])] (table constraint)`

The `UNIQUE` constraint specifies that a group of one or more columns of a table can contain only unique values. The behavior of a unique table constraint is the same as that of a unique column constraint, with the additional capability to span multiple columns. The constraint therefore enforces that any two rows must differ in at least one of these columns.

For the purpose of a unique constraint, null values are not considered equal, unless `NULLS NOT DISTINCT` is specified.

Each unique constraint should name a set of columns that is different from the set of columns named by any other unique or primary key constraint defined for the table. (Otherwise, redundant unique constraints will be discarded.)

When establishing a unique constraint for a multi-level partition hierarchy, all the columns in the partition key of the target partitioned table, as well as those of all its descendant partitioned tables, must be included in the constraint definition.

Adding a unique constraint will automatically create a unique btree index on the column or group of columns used in the constraint. The created index has the same name as the unique constraint.

The optional `INCLUDE` clause adds to that index one or more columns that are simply “payload”: uniqueness is not enforced on them, and the index cannot be searched on the basis of those columns. However they can be retrieved by an index-only scan. Note that although the constraint is not enforced on included columns, it still depends on them. Consequently, some operations on such columns (e.g., `DROP COLUMN`) can cause cascaded constraint and index deletion.

PRIMARY KEY (column constraint)

`PRIMARY KEY (column_name [, ...]) [INCLUDE (column_name [, ...])] (table constraint)`

The `PRIMARY KEY` constraint specifies that a column or columns of a table can contain only unique (non-duplicate), nonnull values. Only one primary key can be specified for a table, whether as a column constraint or a table constraint.

The primary key constraint should name a set of columns that is different from the set of columns named by any unique constraint defined for the same table. (Otherwise, the unique constraint is redundant and will be discarded.)

`PRIMARY KEY` enforces the same data constraints as a combination of `UNIQUE` and `NOT NULL`. However, identifying a set of columns as the primary key also provides metadata about the design of the schema, since a primary key implies that other tables can rely on this set of columns as a unique identifier for rows.

When placed on a partitioned table, `PRIMARY KEY` constraints share the restrictions previously described for `UNIQUE` constraints.

Adding a `PRIMARY KEY` constraint will automatically create a unique btree index on the column or group of columns used in the constraint. That index has the same name as the primary key constraint.

The optional `INCLUDE` clause adds to that index one or more columns that are simply “payload”: uniqueness is not enforced on them, and the index cannot be searched on the basis of those columns. However they can be retrieved by an index-only scan. Note that although the constraint is not enforced on included columns, it still depends on them. Consequently, some operations on such columns (e.g., `DROP COLUMN`) can cause cascaded constraint and index deletion.

`EXCLUDE [USING index_method] (exclude_element WITH operator [, ...]) index_parameters [WHERE (predicate)]`

The `EXCLUDE` clause defines an exclusion constraint, which guarantees that if any two rows are compared on the specified column(s) or expression(s) using the specified operator(s), not all of these comparisons will return `TRUE`. If all of the specified operators test for equality, this is equivalent to a `UNIQUE` constraint, although an ordinary unique constraint will be faster. However, exclusion constraints can specify constraints that are more general than simple equality. For example, you can specify a constraint that no two rows in the table contain overlapping circles (see [Section 8.8](#)) by using the `&&` operator. The operator(s) are required to be commutative.

Exclusion constraints are implemented using an index that has the same name as the constraint, so each specified operator must be associated with an appropriate operator class (see [Section 11.10](#)) for the index access method `index_method`. Each `exclude_element` defines a column of the index, so it can optionally specify a collation, an operator class, operator class parameters, and/or ordering options; these are described fully under [CREATE INDEX](#).

The access method must support `amgettuple` (see [Chapter 63](#)); at present this means GIN cannot be used. Although it's allowed, there is little point in using B-tree or hash indexes with an exclusion constraint, because this does nothing that an ordinary unique constraint doesn't do better. So in practice the access method will always be GiST or SP-GiST.

The `predicate` allows you to specify an exclusion constraint on a subset of the table; internally this creates a partial index. Note that parentheses are required around the predicate.


```
REFERENCES reftable [ ( refcolumn ) ] [ MATCH matchtype ] [ ON DELETE referential_action ] [ ON UPDATE referential_action ] (column constraint)
FOREIGN KEY ( column_name [, ... ] ) REFERENCES reftable [ ( refcolumn [, ... ] ) ] [ MATCH matchtype ] [ ON DELETE referential_action ] [ ON UPDATE referential_action ]
(table constraint)
```

These clauses specify a foreign key constraint, which requires that a group of one or more columns of the new table must only contain values that match values in the referenced column(s) of some row of the referenced table. If the *refcolumn* list is omitted, the primary key of the *reftable* is used. Otherwise, the *refcolumn* list must refer to the columns of a non-deferrable unique or primary key constraint or be the columns of a non-partial unique index. The user must have REFERENCES permission on the referenced table (either the whole table, or the specific referenced columns). The addition of a foreign key constraint requires a SHARE ROW EXCLUSIVE lock on the referenced table. Note that foreign key constraints cannot be defined between temporary tables and permanent tables.

A value inserted into the referencing column(s) is matched against the values of the referenced table and referenced columns using the given match type. There are three match types: MATCH FULL, MATCH PARTIAL, and MATCH SIMPLE (which is the default). MATCH FULL will not allow one column of a multicolumn foreign key to be null unless all foreign key columns are null; if they are all null, the row is not required to have a match in the referenced table. MATCH SIMPLE allows any of the foreign key columns to be null; if any of them are null, the row is not required to have a match in the referenced table. MATCH PARTIAL is not yet implemented. (Of course, NOT NULL constraints can be applied to the referencing column(s) to prevent these cases from arising.)

In addition, when the data in the referenced columns is changed, certain actions are performed on the data in this table's columns. The ON DELETE clause specifies the action to perform when a referenced row in the referenced table is being deleted. Likewise, the ON UPDATE clause specifies the action to perform when a referenced column in the referenced table is being updated to a new value. If the row is updated, but the referenced column is not actually changed, no action is done. Referential actions other than the NO ACTION check cannot be deferred, even if the constraint is declared deferrable. There are the following possible actions for each clause:

NO ACTION

Produce an error indicating that the deletion or update would create a foreign key constraint violation. If the constraint is deferred, this error will be produced at constraint check time if there still exist any referencing rows. This is the default action.

RESTRICT

Produce an error indicating that the deletion or update would create a foreign key constraint violation. This is the same as NO ACTION except that the check is not deferrable.

CASCADE

Delete any rows referencing the deleted row, or update the values of the referencing column(s) to the new values of the referenced columns, respectively.

```
SET NULL [ ( column_name [, ... ] ) ]
```

Set all of the referencing columns, or a specified subset of the referencing columns, to null. A subset of columns can only be specified for ON DELETE actions.

```
SET DEFAULT [ ( column_name [, ... ] ) ]
```

Set all of the referencing columns, or a specified subset of the referencing columns, to their default values. A subset of columns can only be specified for ON DELETE actions. (There must be a row in the referenced table matching the default values, if they are not null, or the operation will fail.)

If the referenced column(s) are changed frequently, it might be wise to add an index to the referencing column(s) so that referential actions associated with the foreign key constraint can be performed more efficiently.

DEFERRABLE

NOT DEFERRABLE

This controls whether the constraint can be deferred. A constraint that is not deferrable will be checked immediately after every command. Checking of constraints that are deferrable can be postponed until the end of the transaction (using the [SET CONSTRAINTS](#) command). `NOT DEFERRABLE` is the default. Currently, only `UNIQUE`, `PRIMARY KEY`, `EXCLUDE`, and `REFERENCES` (foreign key) constraints accept this clause. `NOT NULL` and `CHECK` constraints are not deferrable. Note that deferrable constraints cannot be used as conflict arbitrators in an `INSERT` statement that includes an `ON CONFLICT DO UPDATE` clause.

INITIALLY IMMEDIATE

INITIALLY DEFERRED

If a constraint is deferrable, this clause specifies the default time to check the constraint. If the constraint is `INITIALLY IMMEDIATE`, it is checked after each statement. This is the default. If the constraint is `INITIALLY DEFERRED`, it is checked only at the end of the transaction. The constraint check time can be altered with the [SET CONSTRAINTS](#) command.

USING *method*

This optional clause specifies the table access method to use to store the contents for the new table; the method needs be an access method of type `TABLE`. See [Chapter 62](#) for more information. If this option is not specified, the default table access method is chosen for the new table. See [default_table_access_method](#) for more information.

When creating a partition, the table access method is the access method of its partitioned table, if set.

WITH (*storage_parameter* [= *value*] [, ...])

This clause specifies optional storage parameters for a table or index; see [Storage Parameters](#) below for more information. For backward-compatibility the `WITH` clause for a table can also include `oids=false` to specify that rows of the new table should not contain OIDs (object identifiers), `oids=true` is not supported anymore.

WITHOUT OIDS

This is backward-compatible syntax for declaring a table `WITHOUT OIDS`, creating a table `WITH OIDS` is not supported anymore.

ON COMMIT

The behavior of temporary tables at the end of a transaction block can be controlled using `ON COMMIT`. The three options are:

PRESERVE ROWS

No special action is taken at the ends of transactions. This is the default behavior.

DELETE ROWS

All rows in the temporary table will be deleted at the end of each transaction block. Essentially, an automatic `TRUNCATE` is done at each commit. When used on a partitioned table, this is not cascaded to its partitions.

DROP

The temporary table will be dropped at the end of the current transaction block. When used on a partitioned table, this action drops its partitions and when used on tables with inheritance children, it drops the dependent children.

TABLESPACE *tablespace_name*

The *tablespace_name* is the name of the tablespace in which the new table is to be created. If not specified, [default_tablespace](#) is consulted, or [temp_tablespaces](#) if the table is temporary. For parti-

tioned tables, since no storage is required for the table itself, the tablespace specified overrides `default_tablespace` as the default tablespace to use for any newly created partitions when no other tablespace is explicitly specified.

`USING INDEX TABLESPACE tablespace_name`

This clause allows selection of the tablespace in which the index associated with a `UNIQUE`, `PRIMARY KEY`, or `EXCLUDE` constraint will be created. If not specified, `default_tablespace` is consulted, or `temp_tablespaces` if the table is temporary.

Storage Parameters for Postgres Pro Shardman

Postgres Pro Shardman extends *storage parameters* of tables with its own *storage metaparameters*. They are not stored in the corresponding catalog entry, but are used to tell the Postgres Pro Shardman extension to perform some additional actions. Regular storage parameters are transparently passed to table partitions.

`distributed_by (text)`

This specifies the list of columns to use for the table partitioning. Only hash partitioning is currently supported, so this is effectively an equivalent of `PARTITION BY HASH`, but all the leaf partitions will be created immediately on all replication groups and the table will be registered in the Postgres Pro Shardman metadata.

`num_parts (integer)`

This sets the number of partitions that will be created for this table. This parameter is optional. If it is not specified, for a sharded table, the value of the global setting of `shardman.num_parts` will be used, for a colocated table, the value will be taken from the corresponding colocating table.

`colocate_with (text)`

This specifies the name of the table to colocate with. If set, Postgres Pro Shardman will try to place partitions of the created table with the same partition key on the same nodes as `colocation_table_name`. This parameter is optional.

`partition_by (text)`

This specifies the name of the column to use for the second-level table partitioning. Only range partitioning is currently supported. When this parameter is used, each table partition is created as a partitioned table. Subpartitions can be created immediately if `partition_bounds` parameter is set. This parameter is optional.

`partition_bounds (text)`

This sets bounds of second-level table partitions. Bounds should be a string representation of a two-dimensional array. Each array member is a pair of a lower and upper bound for partitions. If lower and upper bounds are both `NULL`, the default partition is created. Number of partitions is determined by the first array dimension. This parameter is optional.

`global (boolean)`

This defines that the table is global. If set, the table will be distributed on all replication groups and will be synchronized by triggers. This parameter is optional.

Storage Parameters for Not Sharded Tables

The `WITH` clause can specify *storage parameters* for tables, and for indexes associated with a `UNIQUE`, `PRIMARY KEY`, or `EXCLUDE` constraint. Storage parameters for indexes are documented in [CREATE INDEX](#). The storage parameters currently available for tables are listed below. For many of these parameters, as shown, there is an additional parameter with the same name prefixed with `toast.`, which controls the behavior of the table's secondary TOAST table, if any (see [Section 66.2](#) for more information about TOAST). If a table parameter value is set and the equivalent `toast.` parameter is not, the TOAST table

will use the table's parameter value. Specifying these parameters for partitioned tables is not supported, but you may specify them for individual leaf partitions.

`fillfactor (integer)`

The `fillfactor` for a table is a percentage between 10 and 100. 100 (complete packing) is the default. When a smaller `fillfactor` is specified, `INSERT` operations pack table pages only to the indicated percentage; the remaining space on each page is reserved for updating rows on that page. This gives `UPDATE` a chance to place the updated copy of a row on the same page as the original, which is more efficient than placing it on a different page, and makes [heap-only tuple updates](#) more likely. For a table whose entries are never updated, complete packing is the best choice, but in heavily updated tables smaller `fillfactors` are appropriate. This parameter cannot be set for TOAST tables.

`toast_tuple_target (integer)`

The `toast_tuple_target` specifies the minimum tuple length required before we try to compress and/or move long column values into TOAST tables, and is also the target length we try to reduce the length below once toasting begins. This affects columns marked as External (for move), Main (for compression), or Extended (for both) and applies only to new tuples. There is no effect on existing rows. By default this parameter is set to allow at least 4 tuples per block, which with the default block size will be 2040 bytes. Valid values are between 128 bytes and the (block size - header), by default 8160 bytes. Changing this value may not be useful for very short or very long rows. Note that the default setting is often close to optimal, and it is possible that setting this parameter could have negative effects in some cases. This parameter cannot be set for TOAST tables.

`parallel_workers (integer)`

This sets the number of workers that should be used to assist a parallel scan of this table. If not set, the system will determine a value based on the relation size. The actual number of workers chosen by the planner or by utility statements that use parallel scans may be less, for example due to the setting of [max_worker_processes](#).

`autovacuum_enabled, toast.autovacuum_enabled (boolean)`

Enables or disables the autovacuum daemon for a particular table. If true, the autovacuum daemon will perform automatic `VACUUM` and/or `ANALYZE` operations on this table following the rules discussed in [Section 24.1.6](#). If false, this table will not be autovacuumed, except to shrink `pg_xact` and `pg_multixact`. See [Section 24.1.5](#) for more about that. Note that the autovacuum daemon does not run at all (except to shrink `pg_xact` and `pg_multixact`) if the [autovacuum](#) parameter is false; setting individual tables' storage parameters does not override that. Therefore there is seldom much point in explicitly setting this storage parameter to true, only to false.

`vacuum_index_cleanup, toast.vacuum_index_cleanup (enum)`

Forces or disables index cleanup when `VACUUM` is run on this table. The default value is `AUTO`. With `OFF`, index cleanup is disabled, with `ON` it is enabled, and with `AUTO` a decision is made dynamically, each time `VACUUM` runs. The dynamic behavior allows `VACUUM` to avoid needlessly scanning indexes to remove very few dead tuples. Forcibly disabling all index cleanup can speed up `VACUUM` very significantly, but may also lead to severely bloated indexes if table modifications are frequent. The `INDEX_CLEANUP` parameter of [VACUUM](#), if specified, overrides the value of this option.

`vacuum_truncate, toast.vacuum_truncate (boolean)`

Enables or disables vacuum to try to truncate off any empty pages at the end of this table. The default value is `true`. If `true`, `VACUUM` and autovacuum do the truncation and the disk space for the truncated pages is returned to the operating system. Note that the truncation requires `ACCESS EXCLUSIVE` lock on the table. The `TRUNCATE` parameter of [VACUUM](#), if specified, overrides the value of this option.

`parallel_autovacuum_workers (integer)`

The maximum number of additional parallel autovacuum workers that can be launched for table processing. Must be less than or equal to [max_parallel_autovacuum_workers](#). If the value is 0, the

number of table indexes determines the number of parallel autovacuum workers. The default value is -1, which means no parallel vacuuming of the table.

`autovacuum_vacuum_threshold, toast.autovacuum_vacuum_threshold (integer)`

Per-table value for [autovacuum_vacuum_threshold](#) parameter.

`autovacuum_vacuum_scale_factor, toast.autovacuum_vacuum_scale_factor (floating point)`

Per-table value for [autovacuum_vacuum_scale_factor](#) parameter.

`autovacuum_vacuum_insert_threshold, toast.autovacuum_vacuum_insert_threshold (integer)`

Per-table value for [autovacuum_vacuum_insert_threshold](#) parameter. The special value of -1 may be used to disable insert vacuums on the table.

`autovacuum_vacuum_insert_scale_factor, toast.autovacuum_vacuum_insert_scale_factor (floating point)`

Per-table value for [autovacuum_vacuum_insert_scale_factor](#) parameter.

`autovacuum_analyze_threshold (integer)`

Per-table value for [autovacuum_analyze_threshold](#) parameter.

`autovacuum_analyze_scale_factor (floating point)`

Per-table value for [autovacuum_analyze_scale_factor](#) parameter.

`autovacuum_vacuum_cost_delay, toast.autovacuum_vacuum_cost_delay (floating point)`

Per-table value for [autovacuum_vacuum_cost_delay](#) parameter.

`autovacuum_vacuum_cost_limit, toast.autovacuum_vacuum_cost_limit (integer)`

Per-table value for [autovacuum_vacuum_cost_limit](#) parameter.

`autovacuum_freeze_min_age, toast.autovacuum_freeze_min_age (integer)`

Per-table value for [vacuum_freeze_min_age](#) parameter. Note that autovacuum will ignore per-table `autovacuum_freeze_min_age` parameters that are larger than half the system-wide [autovacuum_freeze_max_age](#) setting.

`autovacuum_freeze_max_age, toast.autovacuum_freeze_max_age (integer)`

Per-table value for [autovacuum_freeze_max_age](#) parameter. Note that autovacuum will ignore per-table `autovacuum_freeze_max_age` parameters that are larger than the system-wide setting (it can only be set smaller).

`autovacuum_freeze_table_age, toast.autovacuum_freeze_table_age (integer)`

Per-table value for [vacuum_freeze_table_age](#) parameter.

`autovacuum_multixact_freeze_min_age, toast.autovacuum_multixact_freeze_min_age (integer)`

Per-table value for [vacuum_multixact_freeze_min_age](#) parameter. Note that autovacuum will ignore per-table `autovacuum_multixact_freeze_min_age` parameters that are larger than half the system-wide [autovacuum_multixact_freeze_max_age](#) setting.

`autovacuum_multixact_freeze_max_age, toast.autovacuum_multixact_freeze_max_age (integer)`

Per-table value for [autovacuum_multixact_freeze_max_age](#) parameter. Note that autovacuum will ignore per-table `autovacuum_multixact_freeze_max_age` parameters that are larger than the system-wide setting (it can only be set smaller).

`autovacuum_multixact_freeze_table_age, toast.autovacuum_multixact_freeze_table_age (integer)`

Per-table value for `vacuum_multixact_freeze_table_age` parameter.

`log_autovacuum_min_duration, toast.log_autovacuum_min_duration (integer)`

Per-table value for `log_autovacuum_min_duration` parameter.

`user_catalog_table (boolean)`

Declare the table as an additional catalog table for purposes of logical replication. See [Section 50.6.2](#) for details. This parameter cannot be set for TOAST tables.

Notes

Postgres Pro automatically creates an index for each unique constraint and primary key constraint to enforce uniqueness. Thus, it is not necessary to create an index explicitly for primary key columns. (See [CREATE INDEX](#) for more information.)

Unique constraints and primary keys are not inherited in the current implementation. This makes the combination of inheritance and unique constraints rather dysfunctional.

A table cannot have more than 1600 columns. (In practice, the effective limit is usually lower because of tuple-length constraints.)

Compatibility

The `CREATE TABLE` command conforms to the SQL standard, with exceptions listed below.

Temporary Tables

Note that temporary table usage is highly limited in Postgres Pro Shardman. If a transaction referred a remote sever, it cannot refer a temporary table.

Although the syntax of `CREATE TEMPORARY TABLE` resembles that of the SQL standard, the effect is not the same. In the standard, temporary tables are defined just once and automatically exist (starting with empty contents) in every session that needs them. Postgres Pro instead requires each session to issue its own `CREATE TEMPORARY TABLE` command for each temporary table to be used. This allows different sessions to use the same temporary table name for different purposes, whereas the standard's approach constrains all instances of a given temporary table name to have the same table structure.

The standard's definition of the behavior of temporary tables is widely ignored. Postgres Pro's behavior on this point is similar to that of several other SQL databases.

The SQL standard also distinguishes between global and local temporary tables, where a local temporary table has a separate set of contents for each SQL module within each session, though its definition is still shared across sessions. Since Postgres Pro does not support SQL modules, this distinction is not relevant in Postgres Pro.

For compatibility's sake, Postgres Pro will accept the `GLOBAL` and `LOCAL` keywords in a temporary table declaration, but they currently have no effect. Use of these keywords is discouraged, since future versions of Postgres Pro might adopt a more standard-compliant interpretation of their meaning.

The `ON COMMIT` clause for temporary tables also resembles the SQL standard, but has some differences. If the `ON COMMIT` clause is omitted, SQL specifies that the default behavior is `ON COMMIT DELETE ROWS`. However, the default behavior in Postgres Pro is `ON COMMIT PRESERVE ROWS`. The `ON COMMIT DROP` option does not exist in SQL.

Non-Deferred Uniqueness Constraints

When a `UNIQUE` or `PRIMARY KEY` constraint is not deferrable, Postgres Pro checks for uniqueness immediately whenever a row is inserted or modified. The SQL standard says that uniqueness should be enforced only at the end of the statement; this makes a difference when, for example, a single command updates

multiple key values. To obtain standard-compliant behavior, declare the constraint as `DEFERRABLE` but not deferred (i.e., `INITIALLY IMMEDIATE`). Be aware that this can be significantly slower than immediate uniqueness checking.

Column Check Constraints

The SQL standard says that `CHECK` column constraints can only refer to the column they apply to; only `CHECK` table constraints can refer to multiple columns. Postgres Pro does not enforce this restriction; it treats column and table check constraints alike.

EXCLUDE Constraint

The `EXCLUDE` constraint type is a Postgres Pro extension.

Foreign Key Constraints

The ability to specify column lists in the foreign key actions `SET DEFAULT` and `SET NULL` is a Postgres Pro extension.

It is a Postgres Pro extension that a foreign key constraint may reference columns of a unique index instead of columns of a primary key or unique constraint.

NULL “Constraint”

The `NULL` “constraint” (actually a non-constraint) is a Postgres Pro extension to the SQL standard that is included for compatibility with some other database systems (and for symmetry with the `NOT NULL` constraint). Since it is the default for any column, its presence is simply noise.

Constraint Naming

The SQL standard says that table and domain constraints must have names that are unique across the schema containing the table or domain. Postgres Pro is laxer: it only requires constraint names to be unique across the constraints attached to a particular table or domain. However, this extra freedom does not exist for index-based constraints (`UNIQUE`, `PRIMARY KEY`, and `EXCLUDE` constraints), because the associated index is named the same as the constraint, and index names must be unique across all relations within the same schema.

Currently, Postgres Pro does not record names for not-null constraints at all, so they are not subject to the uniqueness restriction. This might change in a future release.

Inheritance

Multiple inheritance via the `INHERITS` clause is a Postgres Pro language extension. SQL:1999 and later define single inheritance using a different syntax and different semantics. SQL:1999-style inheritance is not yet supported by Postgres Pro.

Zero-Column Tables

Postgres Pro allows a table of no columns to be created (for example, `CREATE TABLE foo();`). This is an extension from the SQL standard, which does not allow zero-column tables. Zero-column tables are not in themselves very useful, but disallowing them creates odd special cases for `ALTER TABLE DROP COLUMN`, so it seems cleaner to ignore this spec restriction.

Multiple Identity Columns

Postgres Pro allows a table to have more than one identity column. The standard specifies that a table can have at most one identity column. This is relaxed mainly to give more flexibility for doing schema changes or migrations. Note that the `INSERT` command supports only one override clause that applies to the entire statement, so having multiple identity columns with different behaviors is not well supported.

Generated Columns

The option `STORED` is not standard but is also used by other SQL implementations. The SQL standard does not specify the storage of generated columns.

LIKE Clause

While a `LIKE` clause exists in the SQL standard, many of the options that Postgres Pro accepts for it are not in the standard, and some of the standard's options are not implemented by Postgres Pro.

WITH Clause

The `WITH` clause is a Postgres Pro extension; storage parameters are not in the standard.

Tablespaces

The Postgres Pro concept of tablespaces is not part of the standard. Hence, the clauses `TABLESPACE` and `USING INDEX TABLESPACE` are extensions.

Typed Tables

Typed tables implement a subset of the SQL standard. According to the standard, a typed table has columns corresponding to the underlying composite type as well as one other column that is the “self-referencing column”. Postgres Pro does not support self-referencing columns explicitly.

PARTITION BY Clause

The `PARTITION BY` clause is a Postgres Pro extension.

PARTITION OF Clause

The `PARTITION OF` clause is a Postgres Pro extension.

CREATE TABLE Limitations for Sharded Tables

The extended `CREATE TABLE` syntax imposes limitations on the general syntax of the command. For example, there is currently no support for:

- Generated columns.
- `REFERENCES` and `FOREIGN KEY` constraints between non-located sharded tables.
- `PARTITION BY` and `PARTITION OF` clauses.
- A foreign key on a global table can reference only another global table and a foreign key on a sharded table can reference a located sharded table or a global table. Note that when a foreign key on a sharded or a global table references a global table, only `NO ACTION` or `RESTRICT` referential actions are supported for the `ON UPDATE` action and only `NO ACTION`, `RESTRICT` or `CASCADE` are supported for the `ON DELETE` action.
- Columns of the `SERIAL8` type are implemented using an automatically created global sequence, so all global sequence properties also apply here.
- For `CREATE TABLE`, all limitations for `ALTER TABLE` apply.
- Using of non-builtin types (types with OIDs ≥ 10000) or non-base types or arrays are not allowed in 'distributed_by' columns.
- Only the global user can create sharded or global table.
- In a located table, the number and types of columns used for table partitioning must be the same as for the table to locate with.
- A temporary table cannot be created as sharded or global.
- Self-referencing sharded tables are allowed only if a foreign key is referencing the same partition of the sharded table.
- Global tables cannot inherit other tables. Other tables should not be inherited from sharded tables.

For tables created using `LIKE source_table` where `source_table` is a local table, the following limitations apply:

- Copying without the *like_option* clause or with `INCLUDING INDEXES` is only supported. Does not apply to global tables.
- With `INCLUDING INDEXES`, only unique indexes and indexes supporting the primary relation key are copied.
- Copying indexes for columns is not supported.
- `EXCLUDE` constraints are not supported.
- Local tables used in `CREATE TABLE LIKE` statement must only have columns of base types.
- Partial indexes are not supported.
- Standard collations are only supported.
- `NULLS NOT DISTINCT` constraint is not supported.

Examples

The following example is for a sharded table. In this example, the table `pgbench_branches` is created, as well as colocated tables `pgbench_accounts` and `pgbench_history`. Each partition of the `pgbench_history` table is additionally subpartitioned by range.

```
CREATE TABLE pgbench_branches (  
    bid integer NOT NULL PRIMARY KEY,  
    bbalance integer,  
    filler character(88)  
)  
WITH (distributed_by = 'bid',  
      num_parts = 8);  
CREATE TABLE pgbench_accounts (  
    aid integer NOT NULL,  
    bid integer,  
    abalance integer,  
    filler character(84),  
    PRIMARY KEY (bid, aid)  
)  
WITH (distributed_by = 'bid',  
      num_parts = 8,  
      colocate_with = 'pgbench_branches');  
CREATE TABLE public.pgbench_history (  
    tid integer,  
    bid integer,  
    aid integer,  
    delta integer,  
    mtime timestamp without time zone,  
    filler character(22)  
)  
WITH (distributed_by = 'bid',  
      colocate_with = 'pgbench_branches',  
      partition_by = 'mtime',  
      partition_bounds =  
          $${{minvalue, '2021-01-01 00:00'}},{ '2021-01-01 00:00', '2022-01-01 00:00'}},  
      {'2022-01-01 00:00', maxvalue}}$$  
);
```

These simple examples of `CREATE TABLE` illustrate [limitations](#) related to creation of colocated tables:

This command creates a table to colocate with:

```
CREATE TABLE teams_players (  
    team_id integer NOT NULL,
```

```
    player_id integer,
    scores int,
    PRIMARY KEY (team_id, player_id)
) WITH (distributed_by='team_id, player_id');
```

This command correctly creates a colocated table:

```
CREATE TABLE players_scores (
    player_id integer NOT NULL,
    team_id integer,
    interval tstzrange,
    scores integer,
    foreign key (team_id, player_id) references teams_players(team_id, player_id)
) WITH (distributed_by='team_id, player_id', colocate_with='teams_players');
```

And this command contains an error in the definition of a foreign key:

```
CREATE TABLE players_scores (
    player_id integer NOT NULL,
    team_id integer,
    interval tstzrange,
    scores integer,
    foreign key (team_id, player_id) references teams_players(team_id, player_id)
) WITH (distributed_by='player_id, team_id', colocate_with='teams_players');
ERROR:  foreign key should start with distributed_by columns
```

Consider another example:

```
CREATE TABLE teams (team_id integer primary key, team_name text) with
(distributed_by='team_id');
CREATE TABLE players_teams (
    player_id integer,
    team_id integer references teams(team_id),
    scores integer
) WITH (distributed_by='player_id', colocate_with='teams');
ERROR:  foreign key should start with distributed_by columns
```

For not sharded tables. Create table films and table distributors:

```
CREATE TABLE films (
    code          char(5) CONSTRAINT firstkey PRIMARY KEY,
    title         varchar(40) NOT NULL,
    did           integer NOT NULL,
    date_prod     date,
    kind          varchar(10),
    len           interval hour to minute
);

CREATE TABLE distributors (
    did          integer PRIMARY KEY GENERATED BY DEFAULT AS IDENTITY,
    name         varchar(40) NOT NULL CHECK (name <> '')
);
```

Create a table with a 2-dimensional array:

```
CREATE TABLE array_int (
    vector int[][]
);
```

Define a unique table constraint for the table films. Unique table constraints can be defined on one or more columns of the table:

```
CREATE TABLE films (
```

```

code          char(5),
title         varchar(40),
did           integer,
date_prod    date,
kind          varchar(10),
len           interval hour to minute,
CONSTRAINT production UNIQUE(date_prod)
);

```

Define a check column constraint:

```

CREATE TABLE distributors (
    did      integer CHECK (did > 100),
    name     varchar(40)
);

```

Define a check table constraint:

```

CREATE TABLE distributors (
    did      integer,
    name     varchar(40),
    CONSTRAINT con1 CHECK (did > 100 AND name <> '')
);

```

Define a primary key table constraint for the table `films`:

```

CREATE TABLE films (
    code      char(5),
    title     varchar(40),
    did       integer,
    date_prod date,
    kind      varchar(10),
    len       interval hour to minute,
    CONSTRAINT code_title PRIMARY KEY(code,title)
);

```

Define a primary key constraint for table `distributors`. The following two examples are equivalent, the first using the table constraint syntax, the second the column constraint syntax:

```

CREATE TABLE distributors (
    did      integer,
    name     varchar(40),
    PRIMARY KEY(did)
);

```

```

CREATE TABLE distributors (
    did      integer PRIMARY KEY,
    name     varchar(40)
);

```

Assign a literal constant default value for the column `name`, arrange for the default value of column `did` to be generated by selecting the next value of a sequence object, and make the default value of `modtime` be the time at which the row is inserted:

```

CREATE TABLE distributors (
    name      varchar(40) DEFAULT 'Luso Films',
    did       integer DEFAULT nextval('distributors_serial'),
    modtime   timestamp DEFAULT current_timestamp
);

```

Define two NOT NULL column constraints on the table `distributors`, one of which is explicitly given a name:

```
CREATE TABLE distributors (  
    did      integer CONSTRAINT no_null NOT NULL,  
    name     varchar(40) NOT NULL  
);
```

Define a unique constraint for the name column:

```
CREATE TABLE distributors (  
    did      integer,  
    name     varchar(40) UNIQUE  
);
```

The same, specified as a table constraint:

```
CREATE TABLE distributors (  
    did      integer,  
    name     varchar(40),  
    UNIQUE(name)  
);
```

Create the same table, specifying 70% fill factor for both the table and its unique index:

```
CREATE TABLE distributors (  
    did      integer,  
    name     varchar(40),  
    UNIQUE(name) WITH (fillfactor=70)  
)  
WITH (fillfactor=70);
```

Create table circles with an exclusion constraint that prevents any two circles from overlapping:

```
CREATE TABLE circles (  
    c circle,  
    EXCLUDE USING gist (c WITH &&)  
);
```

Create table cinemas in tablespace diskvol1:

```
CREATE TABLE cinemas (  
    id serial,  
    name text,  
    location text  
) TABLESPACE diskvol1;
```

Create a composite type and a typed table:

```
CREATE TYPE employee_type AS (name text, salary numeric);  
  
CREATE TABLE employees OF employee_type (  
    PRIMARY KEY (name),  
    salary WITH OPTIONS DEFAULT 1000  
);
```

Create a range partitioned table:

```
CREATE TABLE measurement (  
    logdate      date not null,  
    peaktemp     int,  
    unitsales    int  
) PARTITION BY RANGE (logdate);
```

Create a range partitioned table with multiple columns in the partition key:

```
CREATE TABLE measurement_year_month (  

```

```
    logdate          date not null,  
    peaktemp        int,  
    unitsales       int  
 ) PARTITION BY RANGE (EXTRACT(YEAR FROM logdate), EXTRACT(MONTH FROM logdate));
```

Create a list partitioned table:

```
CREATE TABLE cities (  
    city_id          bigserial not null,  
    name             text not null,  
    population       bigint  
 ) PARTITION BY LIST (left(lower(name), 1));
```

Create a hash partitioned table:

```
CREATE TABLE orders (  
    order_id         bigint not null,  
    cust_id          bigint not null,  
    status           text  
 ) PARTITION BY HASH (order_id);
```

Create partition of a range partitioned table:

```
CREATE TABLE measurement_y2016m07  
    PARTITION OF measurement (  
    unitsales DEFAULT 0  
 ) FOR VALUES FROM ('2016-07-01') TO ('2016-08-01');
```

Create a few partitions of a range partitioned table with multiple columns in the partition key:

```
CREATE TABLE measurement_ym_older  
    PARTITION OF measurement_year_month  
    FOR VALUES FROM (MINVALUE, MINVALUE) TO (2016, 11);  
  
CREATE TABLE measurement_ym_y2016m11  
    PARTITION OF measurement_year_month  
    FOR VALUES FROM (2016, 11) TO (2016, 12);  
  
CREATE TABLE measurement_ym_y2016m12  
    PARTITION OF measurement_year_month  
    FOR VALUES FROM (2016, 12) TO (2017, 01);  
  
CREATE TABLE measurement_ym_y2017m01  
    PARTITION OF measurement_year_month  
    FOR VALUES FROM (2017, 01) TO (2017, 02);
```

Create partition of a list partitioned table:

```
CREATE TABLE cities_ab  
    PARTITION OF cities (  
    CONSTRAINT city_id_nonzero CHECK (city_id != 0)  
 ) FOR VALUES IN ('a', 'b');
```

Create partition of a list partitioned table that is itself further partitioned and then add a partition to it:

```
CREATE TABLE cities_ab  
    PARTITION OF cities (  
    CONSTRAINT city_id_nonzero CHECK (city_id != 0)  
 ) FOR VALUES IN ('a', 'b') PARTITION BY RANGE (population);  
  
CREATE TABLE cities_ab_10000_to_100000  
    PARTITION OF cities_ab FOR VALUES FROM (10000) TO (100000);
```

Create partitions of a hash partitioned table:

```
CREATE TABLE orders_p1 PARTITION OF orders
  FOR VALUES WITH (MODULUS 4, REMAINDER 0);
CREATE TABLE orders_p2 PARTITION OF orders
  FOR VALUES WITH (MODULUS 4, REMAINDER 1);
CREATE TABLE orders_p3 PARTITION OF orders
  FOR VALUES WITH (MODULUS 4, REMAINDER 2);
CREATE TABLE orders_p4 PARTITION OF orders
  FOR VALUES WITH (MODULUS 4, REMAINDER 3);
```

Create a default partition:

```
CREATE TABLE cities_partdef
  PARTITION OF cities DEFAULT;
```

See Also

[ALTER TABLE](#), [DROP TABLE](#), [CREATE TABLE AS](#), [CREATE TABLESPACE](#), [CREATE TYPE](#)

CREATE TABLE AS

CREATE TABLE AS — define a new table from the results of a query

Synopsis

```
CREATE [ [ GLOBAL | LOCAL ] { TEMPORARY | TEMP } | UNLOGGED | CONSTANT ] TABLE [ IF NOT
EXISTS ] table_name
    [ ( column_name [, ...] ) ]
    [ USING method ]
    [ WITH ( storage_parameter [= value] [, ...] ) | WITHOUT OIDS ]
    [ ON COMMIT { PRESERVE ROWS | DELETE ROWS | DROP } ]
    [ TABLESPACE tablespace_name ]
AS query
    [ WITH [ NO ] DATA ]
```

Description

CREATE TABLE AS creates a table and fills it with data computed by a SELECT command. The table columns have the names and data types associated with the output columns of the SELECT (except that you can override the column names by giving an explicit list of new column names).

CREATE TABLE AS bears some resemblance to creating a view, but it is really quite different: it creates a new table and evaluates the query just once to fill the new table initially. The new table will not track subsequent changes to the source tables of the query. In contrast, a view re-evaluates its defining SELECT statement whenever it is queried.

CREATE TABLE AS requires CREATE privilege on the schema used for the table.

Parameters

GLOBAL or LOCAL

Ignored for compatibility. Use of these keywords is deprecated; refer to [CREATE TABLE](#) for details.

TEMPORARY or TEMP

If specified, the table is created as a temporary table. Refer to [CREATE TABLE](#) for details.

UNLOGGED

If specified, the table is created as an unlogged table. Refer to [CREATE TABLE](#) for details.

CONSTANT

If specified, the table is created as read-only. No data can be modified or added to constant tables, and they are not processed by [autovacuum](#). Constant tables cannot be changed to read-write mode.

IF NOT EXISTS

Do not throw an error if a relation with the same name already exists; simply issue a notice and leave the table unmodified.

table_name

The name (optionally schema-qualified) of the table to be created.

column_name

The name of a column in the new table. If column names are not provided, they are taken from the output column names of the query.

USING *method*

This optional clause specifies the table access method to use to store the contents for the new table; the method needs be an access method of type `TABLE`. See [Chapter 62](#) for more information. If this option is not specified, the default table access method is chosen for the new table. See [default_table_access_method](#) for more information.

WITH (*storage_parameter* [= *value*] [, ...])

This clause specifies optional storage parameters for the new table; see [Storage Parameters](#) in the [CREATE TABLE](#) documentation for more information. For backward-compatibility the `WITH` clause for a table can also include `OIDS=FALSE` to specify that rows of the new table should contain no OIDs (object identifiers), `OIDS=TRUE` is not supported anymore.

WITHOUT OIDS

This is backward-compatible syntax for declaring a table `WITHOUT OIDS`, creating a table `WITH OIDS` is not supported anymore.

ON COMMIT

The behavior of temporary tables at the end of a transaction block can be controlled using `ON COMMIT`. The three options are:

PRESERVE ROWS

No special action is taken at the ends of transactions. This is the default behavior.

DELETE ROWS

All rows in the temporary table will be deleted at the end of each transaction block. Essentially, an automatic `TRUNCATE` is done at each commit.

DROP

The temporary table will be dropped at the end of the current transaction block.

TABLESPACE *tablespace_name*

The *tablespace_name* is the name of the tablespace in which the new table is to be created. If not specified, [default_tablespace](#) is consulted, or [temp_tablespaces](#) if the table is temporary.

query

A [SELECT](#), [TABLE](#), or [VALUES](#) command, or an [EXECUTE](#) command that runs a prepared `SELECT`, `TABLE`, or `VALUES` query.

WITH [NO] DATA

This clause specifies whether or not the data produced by the query should be copied into the new table. If not, only the table structure is copied. The default is to copy the data.

Notes

This command is functionally similar to [SELECT INTO](#), but it is preferred since it is less likely to be confused with other uses of the `SELECT INTO` syntax. Furthermore, `CREATE TABLE AS` offers a superset of the functionality offered by `SELECT INTO`.

Examples

Create a new table `films_recent` consisting of only recent entries from the table `films`:

```
CREATE TABLE films_recent AS
SELECT * FROM films WHERE date_prod >= '2002-01-01';
```

To copy a table completely, the short form using the `TABLE` command can also be used:


```
CREATE TABLE films2 AS
TABLE films;
```

Create a new temporary table `films_recent`, consisting of only recent entries from the table `films`, using a prepared statement. The new table will be dropped at commit:

```
PREPARE recentfilms(date) AS
SELECT * FROM films WHERE date_prod > $1;
CREATE TEMP TABLE films_recent ON COMMIT DROP AS
EXECUTE recentfilms('2002-01-01');
```

Compatibility

`CREATE TABLE AS` conforms to the SQL standard. The following are nonstandard extensions:

- The standard requires parentheses around the subquery clause; in Postgres Pro, these parentheses are optional.
- In the standard, the `WITH [NO] DATA` clause is required; in Postgres Pro it is optional.
- Postgres Pro handles temporary tables in a way rather different from the standard; see [CREATE TABLE](#) for details.
- The `WITH` clause is a Postgres Pro extension; storage parameters are not in the standard.
- The Postgres Pro concept of tablespaces is not part of the standard. Hence, the clause `TABLESPACE` is an extension.

See Also

[CREATE MATERIALIZED VIEW](#), [CREATE TABLE](#), [EXECUTE](#), [SELECT](#), [SELECT INTO](#), [VALUES](#)

CREATE TABLESPACE

CREATE TABLESPACE — define a new tablespace

Synopsis

For global tablespaces only:

```
CREATE TABLESPACE tablespace_name
    [ OWNER { new_owner | CURRENT_ROLE | CURRENT_USER | SESSION_USER } ]
    LOCATION 'template'
    [ WITH (global = true [ tablespace_option = value [, ... ] ] ) ]
```

Description

Note that Postgres Pro Shardman only supports global tablespaces.

Postgres Pro Shardman extension of the `CREATE TABLESPACE` syntax enables creation of a new cluster-wide tablespace. All tablespaces in a Postgres Pro Shardman cluster must be cluster-wide. The cluster-wide tablespace is created on each cluster node with the location derived from the `template` parameter.

A tablespace allows database administrators to define an alternative location on the file system where the data files containing database objects (such as tables and indexes) can reside. Only superusers and users with the privileges of the `pg_create_tablespace` role can create tablespaces, but they can assign ownership of tablespaces to non-superusers.

A user with appropriate privileges can pass *tablespace_name* to `CREATE DATABASE`, `CREATE TABLE`, `CREATE INDEX` or `ADD CONSTRAINT` to have the data files for these objects stored within the specified tablespace.

Warning

A tablespace cannot be used independently of the cluster in which it is defined; see [Section 22.7](#).

Parameters

tablespace_name

The name of a tablespace to be created. The name cannot begin with `pg_`, as such names are reserved for system tablespaces.

user_name

The name of the user who will own the tablespace. If omitted, defaults to the user executing the command.

directory

The directory that will be used for the tablespace. The directory must exist (`CREATE TABLESPACE` will not create it), should be empty, and must be owned by the Postgres Pro system user. The directory must be specified by an absolute path name.

global = *true*

Specifies that the created tablespace must be global.

tablespace_option

A tablespace parameter to be set or reset. Currently, the only available parameters are `seq_page_cost`, `random_page_cost`, `effective_io_concurrency`, `maintenance_io_concurrency` and

`compression`. Setting cost value for a particular tablespace will override the planner's usual estimate of the cost of reading pages from tables in that tablespace, and the executor's prefetching behavior, as established by the configuration parameters of the same name (see [seq_page_cost](#), [random_page_cost](#), [effective_io_concurrency](#), [maintenance_io_concurrency](#)). This may be useful if one tablespace is located on a disk which is faster or slower than the remainder of the I/O subsystem. Compression is discussed in section [Chapter 32](#).

Notes

`CREATE TABLESPACE` cannot be executed inside a transaction block.

Examples

To create a tablespace `dbspace` under the file system location `/data/dbs`, first create the directory using operating system facilities on all nodes and set the correct ownership (or ensure that `postgres` user has permissions to create it):

```
mkdir /data/dbs
chown postgres:postgres /data/dbs
```

Then issue the tablespace creation command

```
CREATE TABLESPACE dbspace LOCATION '/data/dbs/ts-{rgid}' WITH (global = true);
```

Compatibility

`CREATE TABLESPACE` is a Postgres Pro extension.

See Also

[CREATE DATABASE](#), [CREATE TABLE](#), [CREATE INDEX](#), [DROP TABLESPACE](#), [ALTER TABLESPACE](#)

CREATE TEXT SEARCH CONFIGURATION

CREATE TEXT SEARCH CONFIGURATION — define a new text search configuration

Synopsis

```
CREATE TEXT SEARCH CONFIGURATION name (  
    PARSER = parser_name |  
    COPY = source_config  
)
```

Description

CREATE TEXT SEARCH CONFIGURATION creates a new text search configuration. A text search configuration specifies a text search parser that can divide a string into tokens, plus dictionaries that can be used to determine which tokens are of interest for searching.

If only the parser is specified, then the new text search configuration initially has no mappings from token types to dictionaries, and therefore will ignore all words. Subsequent ALTER TEXT SEARCH CONFIGURATION commands must be used to create mappings to make the configuration useful. Alternatively, an existing text search configuration can be copied.

If a schema name is given then the text search configuration is created in the specified schema. Otherwise it is created in the current schema.

The user who defines a text search configuration becomes its owner.

Refer to [Chapter 12](#) for further information.

Parameters

name

The name of the text search configuration to be created. The name can be schema-qualified.

parser_name

The name of the text search parser to use for this configuration.

source_config

The name of an existing text search configuration to copy.

Notes

The PARSER and COPY options are mutually exclusive, because when an existing configuration is copied, its parser selection is copied too.

Compatibility

There is no CREATE TEXT SEARCH CONFIGURATION statement in the SQL standard.

See Also

[ALTER TEXT SEARCH CONFIGURATION](#), [DROP TEXT SEARCH CONFIGURATION](#)

CREATE TEXT SEARCH DICTIONARY

CREATE TEXT SEARCH DICTIONARY — define a new text search dictionary

Synopsis

```
CREATE TEXT SEARCH DICTIONARY name (  
    TEMPLATE = template  
    [, option = value [, ... ]]  
)
```

Description

CREATE TEXT SEARCH DICTIONARY creates a new text search dictionary. A text search dictionary specifies a way of recognizing interesting or uninteresting words for searching. A dictionary depends on a text search template, which specifies the functions that actually perform the work. Typically the dictionary provides some options that control the detailed behavior of the template's functions.

If a schema name is given then the text search dictionary is created in the specified schema. Otherwise it is created in the current schema.

The user who defines a text search dictionary becomes its owner.

Refer to [Chapter 12](#) for further information.

Parameters

name

The name of the text search dictionary to be created. The name can be schema-qualified.

template

The name of the text search template that will define the basic behavior of this dictionary.

option

The name of a template-specific option to be set for this dictionary.

value

The value to use for a template-specific option. If the value is not a simple identifier or number, it must be quoted (but you can always quote it, if you wish).

The options can appear in any order.

Examples

The following example command creates a Snowball-based dictionary with a nonstandard list of stop words.

```
CREATE TEXT SEARCH DICTIONARY my_russian (  
    template = snowball,  
    language = russian,  
    stopwords = myrussian  
);
```

Compatibility

There is no CREATE TEXT SEARCH DICTIONARY statement in the SQL standard.

See Also

[ALTER TEXT SEARCH DICTIONARY](#), [DROP TEXT SEARCH DICTIONARY](#)

CREATE TEXT SEARCH PARSER

CREATE TEXT SEARCH PARSER — define a new text search parser

Synopsis

```
CREATE TEXT SEARCH PARSER name (  
    START = start_function ,  
    GETTOKEN = gettoken_function ,  
    END = end_function ,  
    LEXTYPES = lextypes_function  
    [, HEADLINE = headline_function ]  
)
```

Description

CREATE TEXT SEARCH PARSER creates a new text search parser. A text search parser defines a method for splitting a text string into tokens and assigning types (categories) to the tokens. A parser is not particularly useful by itself, but must be bound into a text search configuration along with some text search dictionaries to be used for searching.

If a schema name is given then the text search parser is created in the specified schema. Otherwise it is created in the current schema.

You must be a superuser to use CREATE TEXT SEARCH PARSER. (This restriction is made because an erroneous text search parser definition could confuse or even crash the server.)

Refer to [Chapter 12](#) for further information.

Parameters

name

The name of the text search parser to be created. The name can be schema-qualified.

start_function

The name of the start function for the parser.

gettoken_function

The name of the get-next-token function for the parser.

end_function

The name of the end function for the parser.

lextypes_function

The name of the lextypes function for the parser (a function that returns information about the set of token types it produces).

headline_function

The name of the headline function for the parser (a function that summarizes a set of tokens).

The function names can be schema-qualified if necessary. Argument types are not given, since the argument list for each type of function is predetermined. All except the headline function are required.

The arguments can appear in any order, not only the one shown above.

Compatibility

There is no `CREATE TEXT SEARCH PARSER` statement in the SQL standard.

See Also

[ALTER TEXT SEARCH PARSER](#), [DROP TEXT SEARCH PARSER](#)

CREATE TEXT SEARCH TEMPLATE

CREATE TEXT SEARCH TEMPLATE — define a new text search template

Synopsis

```
CREATE TEXT SEARCH TEMPLATE name (  
    [ INIT = init_function , ]  
    LEXIZE = lexize_function  
)
```

Description

CREATE TEXT SEARCH TEMPLATE creates a new text search template. Text search templates define the functions that implement text search dictionaries. A template is not useful by itself, but must be instantiated as a dictionary to be used. The dictionary typically specifies parameters to be given to the template functions.

If a schema name is given then the text search template is created in the specified schema. Otherwise it is created in the current schema.

You must be a superuser to use CREATE TEXT SEARCH TEMPLATE. This restriction is made because an erroneous text search template definition could confuse or even crash the server. The reason for separating templates from dictionaries is that a template encapsulates the “unsafe” aspects of defining a dictionary. The parameters that can be set when defining a dictionary are safe for unprivileged users to set, and so creating a dictionary need not be a privileged operation.

Refer to [Chapter 12](#) for further information.

Parameters

name

The name of the text search template to be created. The name can be schema-qualified.

init_function

The name of the init function for the template.

lexize_function

The name of the lexize function for the template.

The function names can be schema-qualified if necessary. Argument types are not given, since the argument list for each type of function is predetermined. The lexize function is required, but the init function is optional.

The arguments can appear in any order, not only the one shown above.

Compatibility

There is no CREATE TEXT SEARCH TEMPLATE statement in the SQL standard.

See Also

[ALTER TEXT SEARCH TEMPLATE](#), [DROP TEXT SEARCH TEMPLATE](#)

CREATE TRANSFORM

CREATE TRANSFORM — define a new transform

Synopsis

```
CREATE [ OR REPLACE ] TRANSFORM FOR type_name LANGUAGE lang_name (  
    FROM SQL WITH FUNCTION from_sql_function_name [ (argument_type [, ...]) ],  
    TO SQL WITH FUNCTION to_sql_function_name [ (argument_type [, ...]) ]  
);
```

Description

CREATE TRANSFORM defines a new transform. CREATE OR REPLACE TRANSFORM will either create a new transform, or replace an existing definition.

A transform specifies how to adapt a data type to a procedural language. For example, when writing a function in PL/Python using the `hstore` type, PL/Python has no prior knowledge how to present `hstore` values in the Python environment. Language implementations usually default to using the text representation, but that is inconvenient when, for example, an associative array or a list would be more appropriate.

A transform specifies two functions:

- A “from SQL” function that converts the type from the SQL environment to the language. This function will be invoked on the arguments of a function written in the language.
- A “to SQL” function that converts the type from the language to the SQL environment. This function will be invoked on the return value of a function written in the language.

It is not necessary to provide both of these functions. If one is not specified, the language-specific default behavior will be used if necessary. (To prevent a transformation in a certain direction from happening at all, you could also write a transform function that always errors out.)

To be able to create a transform, you must own and have `USAGE` privilege on the type, have `USAGE` privilege on the language, and own and have `EXECUTE` privilege on the from-SQL and to-SQL functions, if specified.

Parameters

type_name

The name of the data type of the transform.

lang_name

The name of the language of the transform.

from_sql_function_name[(*argument_type* [, ...])]

The name of the function for converting the type from the SQL environment to the language. It must take one argument of type `internal` and return type `internal`. The actual argument will be of the type for the transform, and the function should be coded as if it were. (But it is not allowed to declare an SQL-level function returning `internal` without at least one argument of type `internal`.) The actual return value will be something specific to the language implementation. If no argument list is specified, the function name must be unique in its schema.

to_sql_function_name[(*argument_type* [, ...])]

The name of the function for converting the type from the language to the SQL environment. It must take one argument of type `internal` and return the type that is the type for the transform. The actual argument value will be something specific to the language implementation. If no argument list is specified, the function name must be unique in its schema.

Notes

Use `DROP TRANSFORM` to remove transforms.

Examples

To create a transform for type `hstore` and language `plpython3u`, first set up the type and the language:

```
CREATE TYPE hstore ...;
```

```
CREATE EXTENSION plpython3u;
```

Then create the necessary functions:

```
CREATE FUNCTION hstore_to_plpython(val internal) RETURNS internal
LANGUAGE C STRICT IMMUTABLE
AS ...;
```

```
CREATE FUNCTION plpython_to_hstore(val internal) RETURNS hstore
LANGUAGE C STRICT IMMUTABLE
AS ...;
```

And finally create the transform to connect them all together:

```
CREATE TRANSFORM FOR hstore LANGUAGE plpython3u (
    FROM SQL WITH FUNCTION hstore_to_plpython(internal),
    TO SQL WITH FUNCTION plpython_to_hstore(internal)
);
```

In practice, these commands would be wrapped up in an extension.

The `contrib` section contains a number of extensions that provide transforms, which can serve as real-world examples.

Compatibility

This form of `CREATE TRANSFORM` is a Postgres Pro extension. There is a `CREATE TRANSFORM` command in the SQL standard, but it is for adapting data types to client languages. That usage is not supported by Postgres Pro.

See Also

[CREATE FUNCTION](#), [CREATE LANGUAGE](#), [CREATE TYPE](#), [DROP TRANSFORM](#)

CREATE TRIGGER

CREATE TRIGGER — define a new trigger

Synopsis

```
CREATE [ OR REPLACE ] [ CONSTRAINT ] TRIGGER name { BEFORE | AFTER | INSTEAD OF }
{ event [ OR ... ] }
ON table_name
[ FROM referenced_table_name ]
[ NOT DEFERRABLE | [ DEFERRABLE ] [ INITIALLY IMMEDIATE | INITIALLY DEFERRED ] ]
[ REFERENCING { { OLD | NEW } TABLE [ AS ] transition_relation_name } [ ... ] ]
[ FOR [ EACH ] { ROW | STATEMENT } ]
[ WHEN ( condition ) ]
EXECUTE { FUNCTION | PROCEDURE } function_name ( arguments )
```

where *event* can be one of:

```
INSERT
UPDATE [ OF column_name [, ... ] ]
DELETE
TRUNCATE
```

Description

CREATE TRIGGER creates a new trigger. CREATE OR REPLACE TRIGGER will either create a new trigger, or replace an existing trigger. The trigger will be associated with the specified table, view, or foreign table and will execute the specified function *function_name* when certain operations are performed on that table.

To replace the current definition of an existing trigger, use CREATE OR REPLACE TRIGGER, specifying the existing trigger's name and parent table. All other properties are replaced.

The trigger can be specified to fire before the operation is attempted on a row (before constraints are checked and the INSERT, UPDATE, or DELETE is attempted); or after the operation has completed (after constraints are checked and the INSERT, UPDATE, or DELETE has completed); or instead of the operation (in the case of inserts, updates or deletes on a view). If the trigger fires before or instead of the event, the trigger can skip the operation for the current row, or change the row being inserted (for INSERT and UPDATE operations only). If the trigger fires after the event, all changes, including the effects of other triggers, are “visible” to the trigger.

A trigger that is marked FOR EACH ROW is called once for every row that the operation modifies. For example, a DELETE that affects 10 rows will cause any ON DELETE triggers on the target relation to be called 10 separate times, once for each deleted row. In contrast, a trigger that is marked FOR EACH STATEMENT only executes once for any given operation, regardless of how many rows it modifies (in particular, an operation that modifies zero rows will still result in the execution of any applicable FOR EACH STATEMENT triggers).

Triggers that are specified to fire INSTEAD OF the trigger event must be marked FOR EACH ROW, and can only be defined on views. BEFORE and AFTER triggers on a view must be marked as FOR EACH STATEMENT.

In addition, triggers may be defined to fire for TRUNCATE, though only FOR EACH STATEMENT.

The following table summarizes which types of triggers may be used on tables, views, and foreign tables:

When	Event	Row-level	Statement-level
BEFORE	INSERT/UPDATE/DELETE	Tables and foreign tables	Tables, views, and foreign tables

When	Event	Row-level	Statement-level
	TRUNCATE	—	Tables and foreign tables
AFTER	INSERT/UPDATE/DELETE	Tables and foreign tables	Tables, views, and foreign tables
	TRUNCATE	—	Tables and foreign tables
INSTEAD OF	INSERT/UPDATE/DELETE	Views	—
	TRUNCATE	—	—

Also, a trigger definition can specify a Boolean `WHEN` condition, which will be tested to see whether the trigger should be fired. In row-level triggers the `WHEN` condition can examine the old and/or new values of columns of the row. Statement-level triggers can also have `WHEN` conditions, although the feature is not so useful for them since the condition cannot refer to any values in the table.

If multiple triggers of the same kind are defined for the same event, they will be fired in alphabetical order by name.

When the `CONSTRAINT` option is specified, this command creates a *constraint trigger*. This is the same as a regular trigger except that the timing of the trigger firing can be adjusted using `SET CONSTRAINTS`. Constraint triggers must be `AFTER ROW` triggers on plain tables (not foreign tables). They can be fired either at the end of the statement causing the triggering event, or at the end of the containing transaction; in the latter case they are said to be *deferred*. A pending deferred-trigger firing can also be forced to happen immediately by using `SET CONSTRAINTS`. Constraint triggers are expected to raise an exception when the constraints they implement are violated.

The `REFERENCING` option enables collection of *transition relations*, which are row sets that include all of the rows inserted, deleted, or modified by the current SQL statement. This feature lets the trigger see a global view of what the statement did, not just one row at a time. This option is only allowed for an `AFTER` trigger that is not a constraint trigger; also, if the trigger is an `UPDATE` trigger, it must not specify a *column_name* list. `OLD TABLE` may only be specified once, and only for a trigger that can fire on `UPDATE` or `DELETE`; it creates a transition relation containing the *before-images* of all rows updated or deleted by the statement. Similarly, `NEW TABLE` may only be specified once, and only for a trigger that can fire on `UPDATE` or `INSERT`; it creates a transition relation containing the *after-images* of all rows updated or inserted by the statement.

`SELECT` does not modify any rows so you cannot create `SELECT` triggers. Rules and views may provide workable solutions to problems that seem to need `SELECT` triggers.

Refer to [Chapter 40](#) for more information about triggers.

Parameters

name

The name to give the new trigger. This must be distinct from the name of any other trigger for the same table. The name cannot be schema-qualified — the trigger inherits the schema of its table. For a constraint trigger, this is also the name to use when modifying the trigger's behavior using `SET CONSTRAINTS`.

`BEFORE`

`AFTER`

`INSTEAD OF`

Determines whether the function is called before, after, or instead of the event. A constraint trigger can only be specified as `AFTER`.

event

One of `INSERT`, `UPDATE`, `DELETE`, or `TRUNCATE`; this specifies the event that will fire the trigger. Multiple events can be specified using `OR`, except when transition relations are requested.

For `UPDATE` events, it is possible to specify a list of columns using this syntax:

```
UPDATE OF column_name1 [, column_name2 ... ]
```

The trigger will only fire if at least one of the listed columns is mentioned as a target of the `UPDATE` command or if one of the listed columns is a generated column that depends on a column that is the target of the `UPDATE`.

`INSTEAD OF UPDATE` events do not allow a list of columns. A column list cannot be specified when requesting transition relations, either.

table_name

The name (optionally schema-qualified) of the table, view, or foreign table the trigger is for.

referenced_table_name

The (possibly schema-qualified) name of another table referenced by the constraint. This option is used for foreign-key constraints and is not recommended for general use. This can only be specified for constraint triggers.

`DEFERRABLE`

`NOT DEFERRABLE`

`INITIALLY IMMEDIATE`

`INITIALLY DEFERRED`

The default timing of the trigger. See the [CREATE TABLE](#) documentation for details of these constraint options. This can only be specified for constraint triggers.

`REFERENCING`

This keyword immediately precedes the declaration of one or two relation names that provide access to the transition relations of the triggering statement.

`OLD TABLE`

`NEW TABLE`

This clause indicates whether the following relation name is for the before-image transition relation or the after-image transition relation.

transition_relation_name

The (unqualified) name to be used within the trigger for this transition relation.

`FOR EACH ROW`

`FOR EACH STATEMENT`

This specifies whether the trigger function should be fired once for every row affected by the trigger event, or just once per SQL statement. If neither is specified, `FOR EACH STATEMENT` is the default. Constraint triggers can only be specified `FOR EACH ROW`.

condition

A Boolean expression that determines whether the trigger function will actually be executed. If `WHEN` is specified, the function will only be called if the *condition* returns `true`. In `FOR EACH ROW` triggers, the `WHEN` condition can refer to columns of the old and/or new row values by writing `OLD.column_name` or `NEW.column_name` respectively. Of course, `INSERT` triggers cannot refer to `OLD` and `DELETE` triggers cannot refer to `NEW`.

INSTEAD OF triggers do not support WHEN conditions.

Currently, WHEN expressions cannot contain subqueries.

Note that for constraint triggers, evaluation of the WHEN condition is not deferred, but occurs immediately after the row update operation is performed. If the condition does not evaluate to true then the trigger is not queued for deferred execution.

function_name

A user-supplied function that is declared as taking no arguments and returning type `trigger`, which is executed when the trigger fires.

In the syntax of `CREATE TRIGGER`, the keywords `FUNCTION` and `PROCEDURE` are equivalent, but the referenced function must in any case be a function, not a procedure. The use of the keyword `PROCEDURE` here is historical and deprecated.

arguments

An optional comma-separated list of arguments to be provided to the function when the trigger is executed. The arguments are literal string constants. Simple names and numeric constants can be written here, too, but they will all be converted to strings. Please check the description of the implementation language of the trigger function to find out how these arguments can be accessed within the function; it might be different from normal function arguments.

Notes

To create or replace a trigger on a table, the user must have the `TRIGGER` privilege on the table. The user must also have `EXECUTE` privilege on the trigger function.

Use `DROP TRIGGER` to remove a trigger.

Creating a row-level trigger on a partitioned table will cause an identical “clone” trigger to be created on each of its existing partitions; and any partitions created or attached later will have an identical trigger, too. If there is a conflictly-named trigger on a child partition already, an error occurs unless `CREATE OR REPLACE TRIGGER` is used, in which case that trigger is replaced with a clone trigger. When a partition is detached from its parent, its clone triggers are removed.

A column-specific trigger (one defined using the `UPDATE OF column_name` syntax) will fire when any of its columns are listed as targets in the `UPDATE` command's `SET` list. It is possible for a column's value to change even when the trigger is not fired, because changes made to the row's contents by `BEFORE UPDATE` triggers are not considered. Conversely, a command such as `UPDATE ... SET x = x ...` will fire a trigger on column `x`, even though the column's value did not change.

In a `BEFORE` trigger, the WHEN condition is evaluated just before the function is or would be executed, so using WHEN is not materially different from testing the same condition at the beginning of the trigger function. Note in particular that the `NEW` row seen by the condition is the current value, as possibly modified by earlier triggers. Also, a `BEFORE` trigger's WHEN condition is not allowed to examine the system columns of the `NEW` row (such as `ctid`), because those won't have been set yet.

In an `AFTER` trigger, the WHEN condition is evaluated just after the row update occurs, and it determines whether an event is queued to fire the trigger at the end of statement. So when an `AFTER` trigger's WHEN condition does not return true, it is not necessary to queue an event nor to re-fetch the row at end of statement. This can result in significant speedups in statements that modify many rows, if the trigger only needs to be fired for a few of the rows.

In some cases it is possible for a single SQL command to fire more than one kind of trigger. For instance an `INSERT` with an `ON CONFLICT DO UPDATE` clause may cause both insert and update operations, so it will fire both kinds of triggers as needed. The transition relations supplied to triggers are specific to their event type; thus an `INSERT` trigger will see only the inserted rows, while an `UPDATE` trigger will see only the updated rows.

Row updates or deletions caused by foreign-key enforcement actions, such as `ON UPDATE CASCADE` or `ON DELETE SET NULL`, are treated as part of the SQL command that caused them (note that such actions are never deferred). Relevant triggers on the affected table will be fired, so that this provides another way in which an SQL command might fire triggers not directly matching its type. In simple cases, triggers that request transition relations will see all changes caused in their table by a single original SQL command as a single transition relation. However, there are cases in which the presence of an `AFTER ROW` trigger that requests transition relations will cause the foreign-key enforcement actions triggered by a single SQL command to be split into multiple steps, each with its own transition relation(s). In such cases, any statement-level triggers that are present will be fired once per creation of a transition relation set, ensuring that the triggers see each affected row in a transition relation once and only once.

Statement-level triggers on a view are fired only if the action on the view is handled by a row-level `INSTEAD OF` trigger. If the action is handled by an `INSTEAD` rule, then whatever statements are emitted by the rule are executed in place of the original statement naming the view, so that the triggers that will be fired are those on tables named in the replacement statements. Similarly, if the view is automatically updatable, then the action is handled by automatically rewriting the statement into an action on the view's base table, so that the base table's statement-level triggers are the ones that are fired.

Modifying a partitioned table or a table with inheritance children fires statement-level triggers attached to the explicitly named table, but not statement-level triggers for its partitions or child tables. In contrast, row-level triggers are fired on the rows in affected partitions or child tables, even if they are not explicitly named in the query. If a statement-level trigger has been defined with transition relations named by a `REFERENCING` clause, then before and after images of rows are visible from all affected partitions or child tables. In the case of inheritance children, the row images include only columns that are present in the table that the trigger is attached to.

Currently, row-level triggers with transition relations cannot be defined on partitions or inheritance child tables. Also, triggers on partitioned tables may not be `INSTEAD OF`.

Currently, the `OR REPLACE` option is not supported for constraint triggers.

Replacing an existing trigger within a transaction that has already performed updating actions on the trigger's table is not recommended. Trigger firing decisions, or portions of firing decisions, that have already been made will not be reconsidered, so the effects could be surprising.

There are a few built-in trigger functions that can be used to solve common problems without having to write your own trigger code; see [Section 9.29](#).

Examples

Execute the function `check_account_update` whenever a row of the table `accounts` is about to be updated:

```
CREATE TRIGGER check_update
  BEFORE UPDATE ON accounts
  FOR EACH ROW
  EXECUTE FUNCTION check_account_update();
```

Modify that trigger definition to only execute the function if column `balance` is specified as a target in the `UPDATE` command:

```
CREATE OR REPLACE TRIGGER check_update
  BEFORE UPDATE OF balance ON accounts
  FOR EACH ROW
  EXECUTE FUNCTION check_account_update();
```

This form only executes the function if column `balance` has in fact changed value:

```
CREATE TRIGGER check_update
  BEFORE UPDATE ON accounts
  FOR EACH ROW
```



```
WHEN (OLD.balance IS DISTINCT FROM NEW.balance)
EXECUTE FUNCTION check_account_update();
```

Call a function to log updates of accounts, but only if something changed:

```
CREATE TRIGGER log_update
AFTER UPDATE ON accounts
FOR EACH ROW
WHEN (OLD.* IS DISTINCT FROM NEW.*)
EXECUTE FUNCTION log_account_update();
```

Execute the function `view_insert_row` for each row to insert rows into the tables underlying a view:

```
CREATE TRIGGER view_insert
INSTEAD OF INSERT ON my_view
FOR EACH ROW
EXECUTE FUNCTION view_insert_row();
```

Execute the function `check_transfer_balances_to_zero` for each statement to confirm that the transfer rows offset to a net of zero:

```
CREATE TRIGGER transfer_insert
AFTER INSERT ON transfer
REFERENCING NEW TABLE AS inserted
FOR EACH STATEMENT
EXECUTE FUNCTION check_transfer_balances_to_zero();
```

Execute the function `check_matching_pairs` for each row to confirm that changes are made to matching pairs at the same time (by the same statement):

```
CREATE TRIGGER paired_items_update
AFTER UPDATE ON paired_items
REFERENCING NEW TABLE AS newtab OLD TABLE AS oldtab
FOR EACH ROW
EXECUTE FUNCTION check_matching_pairs();
```

[Section 40.4](#) contains a complete example of a trigger function written in C.

Compatibility

The `CREATE TRIGGER` statement in Postgres Pro implements a subset of the SQL standard. The following functionalities are currently missing:

- While transition table names for `AFTER` triggers are specified using the `REFERENCING` clause in the standard way, the row variables used in `FOR EACH ROW` triggers may not be specified in a `REFERENCING` clause. They are available in a manner that is dependent on the language in which the trigger function is written, but is fixed for any one language. Some languages effectively behave as though there is a `REFERENCING` clause containing `OLD ROW AS OLD NEW ROW AS NEW`.
- The standard allows transition tables to be used with column-specific `UPDATE` triggers, but then the set of rows that should be visible in the transition tables depends on the trigger's column list. This is not currently implemented by Postgres Pro.
- Postgres Pro only allows the execution of a user-defined function for the triggered action. The standard allows the execution of a number of other SQL commands, such as `CREATE TABLE`, as the triggered action. This limitation is not hard to work around by creating a user-defined function that executes the desired commands.

SQL specifies that multiple triggers should be fired in time-of-creation order. Postgres Pro uses name order, which was judged to be more convenient.

SQL specifies that `BEFORE DELETE` triggers on cascaded deletes fire *after* the cascaded `DELETE` completes. The Postgres Pro behavior is for `BEFORE DELETE` to always fire before the delete action, even a cascading one. This is considered more consistent. There is also nonstandard behavior if `BEFORE` triggers

modify rows or prevent updates during an update that is caused by a referential action. This can lead to constraint violations or stored data that does not honor the referential constraint.

The ability to specify multiple actions for a single trigger using `OR` is a Postgres Pro extension of the SQL standard.

The ability to fire triggers for `TRUNCATE` is a Postgres Pro extension of the SQL standard, as is the ability to define statement-level triggers on views.

`CREATE CONSTRAINT TRIGGER` is a Postgres Pro extension of the SQL standard. So is the `OR REPLACE` option.

See Also

[ALTER TRIGGER](#), [DROP TRIGGER](#), [CREATE FUNCTION](#), [SET CONSTRAINTS](#)

CREATE TYPE

CREATE TYPE — define a new data type

Synopsis

```
CREATE TYPE name AS
    ( [ attribute_name data_type [ COLLATE collation ] [, ... ] ] )
```

```
CREATE TYPE name AS ENUM
    ( [ 'label' [, ... ] ] )
```

```
CREATE TYPE name AS RANGE (
    SUBTYPE = subtype
    [ , SUBTYPE_OPCLASS = subtype_operator_class ]
    [ , COLLATION = collation ]
    [ , CANONICAL = canonical_function ]
    [ , SUBTYPE_DIFF = subtype_diff_function ]
    [ , MULTIRANGE_TYPE_NAME = multirange_type_name ]
)
```

```
CREATE TYPE name (
    INPUT = input_function,
    OUTPUT = output_function
    [ , RECEIVE = receive_function ]
    [ , SEND = send_function ]
    [ , TYPMOD_IN = type_modifier_input_function ]
    [ , TYPMOD_OUT = type_modifier_output_function ]
    [ , ANALYZE = analyze_function ]
    [ , SUBSCRIPT = subscript_function ]
    [ , INTERNALLENGTH = { internallength | VARIABLE } ]
    [ , PASSEDBYVALUE ]
    [ , ALIGNMENT = alignment ]
    [ , STORAGE = storage ]
    [ , LIKE = like_type ]
    [ , CATEGORY = category ]
    [ , PREFERRED = preferred ]
    [ , DEFAULT = default ]
    [ , ELEMENT = element ]
    [ , DELIMITER = delimiter ]
    [ , COLLATABLE = collatable ]
)
```

```
CREATE TYPE name
```

Description

CREATE TYPE registers a new data type for use in the current database. The user who defines a type becomes its owner.

If a schema name is given then the type is created in the specified schema. Otherwise it is created in the current schema. The type name must be distinct from the name of any existing type or domain in the same schema. (Because tables have associated data types, the type name must also be distinct from the name of any existing table in the same schema.)

There are five forms of CREATE TYPE, as shown in the syntax synopsis above. They respectively create a *composite type*, an *enum type*, a *range type*, a *base type*, or a *shell type*. The first four of these are

discussed in turn below. A shell type is simply a placeholder for a type to be defined later; it is created by issuing `CREATE TYPE` with no parameters except for the type name. Shell types are needed as forward references when creating range types and base types, as discussed in those sections.

Composite Types

The first form of `CREATE TYPE` creates a composite type. The composite type is specified by a list of attribute names and data types. An attribute's collation can be specified too, if its data type is collatable. A composite type is essentially the same as the row type of a table, but using `CREATE TYPE` avoids the need to create an actual table when all that is wanted is to define a type. A stand-alone composite type is useful, for example, as the argument or return type of a function.

To be able to create a composite type, you must have `USAGE` privilege on all attribute types.

Enumerated Types

The second form of `CREATE TYPE` creates an enumerated (enum) type, as described in [Section 8.7](#). Enum types take a list of quoted labels, each of which must be less than `NAMEDATALEN` bytes long (64 bytes in a standard Postgres Pro build). (It is possible to create an enumerated type with zero labels, but such a type cannot be used to hold values before at least one label is added using `ALTER TYPE`.)

Range Types

The third form of `CREATE TYPE` creates a new range type, as described in [Section 8.17](#).

The range type's *subtype* can be any type with an associated b-tree operator class (to determine the ordering of values for the range type). Normally the subtype's default b-tree operator class is used to determine ordering; to use a non-default operator class, specify its name with *subtype_opclass*. If the subtype is collatable, and you want to use a non-default collation in the range's ordering, specify the desired collation with the *collation* option.

The optional *canonical* function must take one argument of the range type being defined, and return a value of the same type. This is used to convert range values to a canonical form, when applicable. See [Section 8.17.8](#) for more information. Creating a *canonical* function is a bit tricky, since it must be defined before the range type can be declared. To do this, you must first create a shell type, which is a placeholder type that has no properties except a name and an owner. This is done by issuing the command `CREATE TYPE name`, with no additional parameters. Then the function can be declared using the shell type as argument and result, and finally the range type can be declared using the same name. This automatically replaces the shell type entry with a valid range type.

The optional *subtype_diff* function must take two values of the *subtype* type as argument, and return a double precision value representing the difference between the two given values. While this is optional, providing it allows much greater efficiency of GiST indexes on columns of the range type. See [Section 8.17.8](#) for more information.

The optional *multirange_type_name* parameter specifies the name of the corresponding multirange type. If not specified, this name is chosen automatically as follows. If the range type name contains the substring *range*, then the multirange type name is formed by replacement of the *range* substring with *multirange* in the range type name. Otherwise, the multirange type name is formed by appending a *_multirange* suffix to the range type name.

Base Types

The fourth form of `CREATE TYPE` creates a new base type (scalar type). To create a new base type, you must be a superuser. (This restriction is made because an erroneous type definition could confuse or even crash the server.)

The parameters can appear in any order, not only that illustrated above, and most are optional. You must register two or more functions (using `CREATE FUNCTION`) before defining the type. The support functions *input_function* and *output_function* are required, while the functions *receive_function*,

send_function, *type_modifier_input_function*, *type_modifier_output_function*, *analyze_function*, and *subscript_function* are optional. Generally these functions have to be coded in C or another low-level language.

The *input_function* converts the type's external textual representation to the internal representation used by the operators and functions defined for the type. *output_function* performs the reverse transformation. The input function can be declared as taking one argument of type *cstring*, or as taking three arguments of types *cstring*, *oid*, *integer*. The first argument is the input text as a C string, the second argument is the type's own OID (except for array types, which instead receive their element type's OID), and the third is the *typmod* of the destination column, if known (-1 will be passed if not). The input function must return a value of the data type itself. Usually, an input function should be declared *STRICT*; if it is not, it will be called with a NULL first parameter when reading a NULL input value. The function must still return NULL in this case, unless it raises an error. (This case is mainly meant to support domain input functions, which might need to reject NULL inputs.) The output function must be declared as taking one argument of the new data type. The output function must return type *cstring*. Output functions are not invoked for NULL values.

The optional *receive_function* converts the type's external binary representation to the internal representation. If this function is not supplied, the type cannot participate in binary input. The binary representation should be chosen to be cheap to convert to internal form, while being reasonably portable. (For example, the standard integer data types use network byte order as the external binary representation, while the internal representation is in the machine's native byte order.) The receive function should perform adequate checking to ensure that the value is valid. The receive function can be declared as taking one argument of type *internal*, or as taking three arguments of types *internal*, *oid*, *integer*. The first argument is a pointer to a *StringInfo* buffer holding the received byte string; the optional arguments are the same as for the text input function. The receive function must return a value of the data type itself. Usually, a receive function should be declared *STRICT*; if it is not, it will be called with a NULL first parameter when reading a NULL input value. The function must still return NULL in this case, unless it raises an error. (This case is mainly meant to support domain receive functions, which might need to reject NULL inputs.) Similarly, the optional *send_function* converts from the internal representation to the external binary representation. If this function is not supplied, the type cannot participate in binary output. The send function must be declared as taking one argument of the new data type. The send function must return type *bytea*. Send functions are not invoked for NULL values.

You should at this point be wondering how the input and output functions can be declared to have results or arguments of the new type, when they have to be created before the new type can be created. The answer is that the type should first be defined as a *shell type*, which is a placeholder type that has no properties except a name and an owner. This is done by issuing the command `CREATE TYPE name`, with no additional parameters. Then the C I/O functions can be defined referencing the shell type. Finally, `CREATE TYPE` with a full definition replaces the shell entry with a complete, valid type definition, after which the new type can be used normally.

The optional *type_modifier_input_function* and *type_modifier_output_function* are needed if the type supports modifiers, that is optional constraints attached to a type declaration, such as `char(5)` or `numeric(30,2)`. Postgres Pro allows user-defined types to take one or more simple constants or identifiers as modifiers. However, this information must be capable of being packed into a single non-negative integer value for storage in the system catalogs. The *type_modifier_input_function* is passed the declared modifier(s) in the form of a *cstring* array. It must check the values for validity (throwing an error if they are wrong), and if they are correct, return a single non-negative integer value that will be stored as the column "typmod". Type modifiers will be rejected if the type does not have a *type_modifier_input_function*. The *type_modifier_output_function* converts the internal integer typmod value back to the correct form for user display. It must return a *cstring* value that is the exact string to append to the type name; for example `numeric`'s function might return `(30,2)`. It is allowed to omit the *type_modifier_output_function*, in which case the default display format is just the stored typmod integer value enclosed in parentheses.

The optional *analyze_function* performs type-specific statistics collection for columns of the data type. By default, `ANALYZE` will attempt to gather statistics using the type's "equals" and "less-than" operators,

if there is a default b-tree operator class for the type. For non-scalar types this behavior is likely to be unsuitable, so it can be overridden by specifying a custom analysis function. The analysis function must be declared to take a single argument of type `internal`, and return a boolean result.

The optional `subscript_function` allows the data type to be subscripted in SQL commands. Specifying this function does not cause the type to be considered a “true” array type; for example, it will not be a candidate for the result type of `ARRAY[]` constructs. But if subscripting a value of the type is a natural notation for extracting data from it, then a `subscript_function` can be written to define what that means. The subscript function must be declared to take a single argument of type `internal`, and return an `internal` result, which is a pointer to a struct of methods (functions) that implement subscripting. The detailed API for subscript functions appears in [src/include/nodes/subscripting.h](#). Additional information appears in [Array Types](#) below.

While the details of the new type's internal representation are only known to the I/O functions and other functions you create to work with the type, there are several properties of the internal representation that must be declared to Postgres Pro. Foremost of these is `internallength`. Base data types can be fixed-length, in which case `internallength` is a positive integer, or variable-length, indicated by setting `internallength` to `VARIABLE`. (Internally, this is represented by setting `typlen` to -1.) The internal representation of all variable-length types must start with a 4-byte integer giving the total length of this value of the type. (Note that the length field is often encoded, as described in [Section 66.2](#); it's unwise to access it directly.)

The optional flag `PASSEDBYVALUE` indicates that values of this data type are passed by value, rather than by reference. Types passed by value must be fixed-length, and their internal representation cannot be larger than the size of the `Datum` type (4 bytes on some machines, 8 bytes on others).

The `alignment` parameter specifies the storage alignment required for the data type. The allowed values equate to alignment on 1, 2, 4, or 8 byte boundaries. Note that variable-length types must have an alignment of at least 4, since they necessarily contain an `int4` as their first component.

The `storage` parameter allows selection of storage strategies for variable-length data types. (Only `plain` is allowed for fixed-length types.) `plain` specifies that data of the type will always be stored in-line and not compressed. `extended` specifies that the system will first try to compress a long data value, and will move the value out of the main table row if it's still too long. `external` allows the value to be moved out of the main table, but the system will not try to compress it. `main` allows compression, but discourages moving the value out of the main table. (Data items with this storage strategy might still be moved out of the main table if there is no other way to make a row fit, but they will be kept in the main table preferentially over `extended` and `external` items.)

All `storage` values other than `plain` imply that the functions of the data type can handle values that have been *toasted*, as described in [Section 66.2](#) and [Section 39.13.1](#). The specific other value given merely determines the default TOAST storage strategy for columns of a toastable data type; users can pick other strategies for individual columns using `ALTER TABLE SET STORAGE`.

The `like_type` parameter provides an alternative method for specifying the basic representation properties of a data type: copy them from some existing type. The values of `internallength`, `passedbyvalue`, `alignment`, and `storage` are copied from the named type. (It is possible, though usually undesirable, to override some of these values by specifying them along with the `LIKE` clause.) Specifying representation this way is especially useful when the low-level implementation of the new type “piggybacks” on an existing type in some fashion.

The `category` and `preferred` parameters can be used to help control which implicit cast will be applied in ambiguous situations. Each data type belongs to a category named by a single ASCII character, and each type is either “preferred” or not within its category. The parser will prefer casting to preferred types (but only from other types within the same category) when this rule is helpful in resolving overloaded functions or operators. For more details see [Chapter 10](#). For types that have no implicit casts to or from any other types, it is sufficient to leave these settings at the defaults. However, for a group of related types that have implicit casts, it is often helpful to mark them all as belonging to a category and select one or two of the “most general” types as being preferred within the category. The `category` parameter

is especially useful when adding a user-defined type to an existing built-in category, such as the numeric or string types. However, it is also possible to create new entirely-user-defined type categories. Select any ASCII character other than an upper-case letter to name such a category.

A default value can be specified, in case a user wants columns of the data type to default to something other than the null value. Specify the default with the `DEFAULT` key word. (Such a default can be overridden by an explicit `DEFAULT` clause attached to a particular column.)

To indicate that a type is a fixed-length array type, specify the type of the array elements using the `ELEMENT` key word. For example, to define an array of 4-byte integers (`int4`), specify `ELEMENT = int4`. For more details, see [Array Types](#) below.

To indicate the delimiter to be used between values in the external representation of arrays of this type, *delimiter* can be set to a specific character. The default delimiter is the comma (`,`). Note that the delimiter is associated with the array element type, not the array type itself.

If the optional Boolean parameter *collatable* is true, column definitions and expressions of the type may carry collation information through use of the `COLLATE` clause. It is up to the implementations of the functions operating on the type to actually make use of the collation information; this does not happen automatically merely by marking the type collatable.

Array Types

Whenever a user-defined type is created, Postgres Pro automatically creates an associated array type, whose name consists of the element type's name prepended with an underscore, and truncated if necessary to keep it less than `NAMEDATALEN` bytes long. (If the name so generated collides with an existing type name, the process is repeated until a non-colliding name is found.) This implicitly-created array type is variable length and uses the built-in input and output functions `array_in` and `array_out`. Furthermore, this type is what the system uses for constructs such as `ARRAY[]` over the user-defined type. The array type tracks any changes in its element type's owner or schema, and is dropped if the element type is.

You might reasonably ask why there is an `ELEMENT` option, if the system makes the correct array type automatically. The main case where it's useful to use `ELEMENT` is when you are making a fixed-length type that happens to be internally an array of a number of identical things, and you want to allow these things to be accessed directly by subscripting, in addition to whatever operations you plan to provide for the type as a whole. For example, type `point` is represented as just two floating-point numbers, which can be accessed using `point[0]` and `point[1]`. Note that this facility only works for fixed-length types whose internal form is exactly a sequence of identical fixed-length fields. For historical reasons (i.e., this is clearly wrong but it's far too late to change it), subscripting of fixed-length array types starts from zero, rather than from one as for variable-length arrays.

Specifying the `SUBSCRIPT` option allows a data type to be subscripted, even though the system does not otherwise regard it as an array type. The behavior just described for fixed-length arrays is actually implemented by the `SUBSCRIPT` handler function `raw_array_subscript_handler`, which is used automatically if you specify `ELEMENT` for a fixed-length type without also writing `SUBSCRIPT`.

When specifying a custom `SUBSCRIPT` function, it is not necessary to specify `ELEMENT` unless the `SUBSCRIPT` handler function needs to consult `typelem` to find out what to return. Be aware that specifying `ELEMENT` causes the system to assume that the new type contains, or is somehow physically dependent on, the element type; thus for example changing properties of the element type won't be allowed if there are any columns of the dependent type.

Parameters

name

The name (optionally schema-qualified) of a type to be created.

attribute_name

The name of an attribute (column) for the composite type.

data_type

The name of an existing data type to become a column of the composite type.

collation

The name of an existing collation to be associated with a column of a composite type, or with a range type.

label

A string literal representing the textual label associated with one value of an enum type.

subtype

The name of the element type that the range type will represent ranges of.

subtype_operator_class

The name of a b-tree operator class for the subtype.

canonical_function

The name of the canonicalization function for the range type.

subtype_diff_function

The name of a difference function for the subtype.

multirange_type_name

The name of the corresponding multirange type.

input_function

The name of a function that converts data from the type's external textual form to its internal form.

output_function

The name of a function that converts data from the type's internal form to its external textual form.

receive_function

The name of a function that converts data from the type's external binary form to its internal form.

send_function

The name of a function that converts data from the type's internal form to its external binary form.

type_modifier_input_function

The name of a function that converts an array of modifier(s) for the type into internal form.

type_modifier_output_function

The name of a function that converts the internal form of the type's modifier(s) to external textual form.

analyze_function

The name of a function that performs statistical analysis for the data type.

subscript_function

The name of a function that defines what subscripting a value of the data type does.

internallength

A numeric constant that specifies the length in bytes of the new type's internal representation. The default assumption is that it is variable-length.

alignment

The storage alignment requirement of the data type. If specified, it must be `char`, `int2`, `int4`, or `double`; the default is `int4`.

storage

The storage strategy for the data type. If specified, must be `plain`, `external`, `extended`, or `main`; the default is `plain`.

like_type

The name of an existing data type that the new type will have the same representation as. The values of *internallength*, *passedbyvalue*, *alignment*, and *storage* are copied from that type, unless overridden by explicit specification elsewhere in this `CREATE TYPE` command.

category

The category code (a single ASCII character) for this type. The default is `'U'` for “user-defined type”. Other standard category codes can be found in [Table 54.67](#). You may also choose other ASCII characters in order to create custom categories.

preferred

True if this type is a preferred type within its type category, else false. The default is false. Be very careful about creating a new preferred type within an existing type category, as this could cause surprising changes in behavior.

default

The default value for the data type. If this is omitted, the default is null.

element

The type being created is an array; this specifies the type of the array elements.

delimiter

The delimiter character to be used between values in arrays made of this type.

collatable

True if this type's operations can use collation information. The default is false.

Notes

Because there are no restrictions on use of a data type once it's been created, creating a base type or range type is tantamount to granting public execute permission on the functions mentioned in the type definition. This is usually not an issue for the sorts of functions that are useful in a type definition. But you might want to think twice before designing a type in a way that would require “secret” information to be used while converting it to or from external form.

Before PostgreSQL version 8.3, the name of a generated array type was always exactly the element type's name with one underscore character (`_`) prepended. (Type names were therefore restricted in length to one fewer character than other names.) While this is still usually the case, the array type name may vary from this in case of maximum-length names or collisions with user type names that begin with underscore. Writing code that depends on this convention is therefore deprecated. Instead, use `pg_type.typarray` to locate the array type associated with a given type.

It may be advisable to avoid using type and table names that begin with underscore. While the server will change generated array type names to avoid collisions with user-given names, there is still risk of confusion, particularly with old client software that may assume that type names beginning with underscores always represent arrays.

Before PostgreSQL version 8.2, the shell-type creation syntax `CREATE TYPE name` did not exist. The way to create a new base type was to create its input function first. In this approach, Postgres Pro will first see the name of the new data type as the return type of the input function. The shell type is implicitly created in this situation, and then it can be referenced in the definitions of the remaining I/O functions. This approach still works, but is deprecated and might be disallowed in some future release. Also, to avoid accidentally cluttering the catalogs with shell types as a result of simple typos in function definitions, a shell type will only be made this way when the input function is written in C.

In Postgres Pro version 16 and later, it is desirable for base types' input functions to return “soft” errors using the new `errsave()/ereturn()` mechanism, rather than throwing `ereport()` exceptions as in previous versions.

Examples

This example creates a composite type and uses it in a function definition:

```
CREATE TYPE compfoo AS (f1 int, f2 text);

CREATE FUNCTION getfoo() RETURNS SETOF compfoo AS $$
    SELECT fooid, foename FROM foo
$$ LANGUAGE SQL;
```

This example creates an enumerated type and uses it in a table definition:

```
CREATE TYPE bug_status AS ENUM ('new', 'open', 'closed');

CREATE TABLE bug (
    id serial,
    description text,
    status bug_status
);
```

This example creates a range type:

```
CREATE TYPE float8_range AS RANGE (subtype = float8, subtype_diff = float8mi);
```

This example creates the base data type `box` and then uses the type in a table definition:

```
CREATE TYPE box;

CREATE FUNCTION my_box_in_function(cstring) RETURNS box AS ... ;
CREATE FUNCTION my_box_out_function(box) RETURNS cstring AS ... ;

CREATE TYPE box (
    INTERNALLENGTH = 16,
    INPUT = my_box_in_function,
    OUTPUT = my_box_out_function
);

CREATE TABLE myboxes (
    id integer,
    description box
);
```

If the internal structure of `box` were an array of four `float4` elements, we might instead use:

```
CREATE TYPE box (
    INTERNALLENGTH = 16,
    INPUT = my_box_in_function,
    OUTPUT = my_box_out_function,
    ELEMENT = float4
);
```

which would allow a box value's component numbers to be accessed by subscripting. Otherwise the type behaves the same as before.

This example creates a large object type and uses it in a table definition:

```
CREATE TYPE bigobj (  
    INPUT = lo_filein, OUTPUT = lo_fileout,  
    INTERNALLENGTH = VARIABLE  
);  
CREATE TABLE big_objs (  
    id integer,  
    obj bigobj  
);
```

More examples, including suitable input and output functions, are in [Section 39.13](#).

Compatibility

The first form of the `CREATE TYPE` command, which creates a composite type, conforms to the SQL standard. The other forms are Postgres Pro extensions. The `CREATE TYPE` statement in the SQL standard also defines other forms that are not implemented in Postgres Pro.

The ability to create a composite type with zero attributes is a Postgres Pro-specific deviation from the standard (analogous to the same case in `CREATE TABLE`).

See Also

[ALTER TYPE](#), [CREATE DOMAIN](#), [CREATE FUNCTION](#), [DROP TYPE](#)

CREATE USER

CREATE USER — define a new database role

Synopsis

```
CREATE USER name [ [ WITH ] option [ ... ] ]
```

where *option* can be:

```
SUPERUSER | NOSUPERUSER
| CREATEDB | NOCREATEDB
| CREATEROLE | NOCREATEROLE
| INHERIT | NOINHERIT
| LOGIN | NOLOGIN
| REPLICATION | NOREPLICATION
| BYPASSRLS | NOBYPASSRLS
| CONNECTION LIMIT connlimit
| [ ENCRYPTED ] PASSWORD 'password' | PASSWORD NULL
| VALID UNTIL 'timestamp'
| IN ROLE role_name [, ...]
| IN GROUP role_name [, ...]
| ROLE role_name [, ...]
| ADMIN role_name [, ...]
| USER role_name [, ...]
| SYSID uid
| PROFILE profile_name
```

Description

CREATE USER is now an alias for [CREATE ROLE](#). The only difference is that when the command is spelled CREATE USER, LOGIN is assumed by default, whereas NOLOGIN is assumed when the command is spelled CREATE ROLE.

Compatibility

The CREATE USER statement is a Postgres Pro extension. The SQL standard leaves the definition of users to the implementation.

See Also

[CREATE ROLE](#)

CREATE VIEW

CREATE VIEW — define a new view

Synopsis

```
CREATE [ OR REPLACE ] [ TEMP | TEMPORARY ] [ RECURSIVE ] VIEW name [ ( column_name
[ , ... ] ) ]
    [ WITH ( view_option_name [= view_option_value] [ , ... ] ) ]
    AS query
    [ WITH [ CASCADED | LOCAL ] CHECK OPTION ]
```

Description

CREATE VIEW defines a view of a query. The view is not physically materialized. Instead, the query is run every time the view is referenced in a query.

CREATE OR REPLACE VIEW is similar, but if a view of the same name already exists, it is replaced. The new query must generate the same columns that were generated by the existing view query (that is, the same column names in the same order and with the same data types), but it may add additional columns to the end of the list. The calculations giving rise to the output columns may be completely different.

If a schema name is given (for example, CREATE VIEW myschema.myview ...) then the view is created in the specified schema. Otherwise it is created in the current schema. Temporary views exist in a special schema, so a schema name cannot be given when creating a temporary view. The name of the view must be distinct from the name of any other relation (table, sequence, index, view, materialized view, or foreign table) in the same schema.

Parameters

TEMPORARY or TEMP

If specified, the view is created as a temporary view. Temporary views are automatically dropped at the end of the current session. Existing permanent relations with the same name are not visible to the current session while the temporary view exists, unless they are referenced with schema-qualified names.

If any of the tables referenced by the view are temporary, the view is created as a temporary view (whether TEMPORARY is specified or not).

RECURSIVE

Creates a recursive view. The syntax

```
CREATE RECURSIVE VIEW [ schema . ] view_name (column_names) AS SELECT ...;
```

is equivalent to

```
CREATE VIEW [ schema . ] view_name AS WITH RECURSIVE view_name (column_names) AS
    (SELECT ...) SELECT column_names FROM view_name;
```

A view column name list must be specified for a recursive view.

name

The name (optionally schema-qualified) of a view to be created.

column_name

An optional list of names to be used for columns of the view. If not given, the column names are deduced from the query.

```
WITH ( view_option_name [= view_option_value] [, ... ] )
```

This clause specifies optional parameters for a view; the following parameters are supported:

`check_option` (enum)

This parameter may be either `local` or `cascaded`, and is equivalent to specifying `WITH [CASCADED | LOCAL] CHECK OPTION` (see below).

`security_barrier` (boolean)

This should be used if the view is intended to provide row-level security. See [Section 42.5](#) for full details.

`security_invoker` (boolean)

This option causes the underlying base relations to be checked against the privileges of the user of the view rather than the view owner. See the notes below for full details.

All of the above options can be changed on existing views using [ALTER VIEW](#).

query

A [SELECT](#) or [VALUES](#) command which will provide the columns and rows of the view.

```
WITH [ CASCADED | LOCAL ] CHECK OPTION
```

This option controls the behavior of automatically updatable views. When this option is specified, `INSERT`, `UPDATE`, and `MERGE` commands on the view will be checked to ensure that new rows satisfy the view-defining condition (that is, the new rows are checked to ensure that they are visible through the view). If they are not, the update will be rejected. If the `CHECK OPTION` is not specified, `INSERT`, `UPDATE`, and `MERGE` commands on the view are allowed to create rows that are not visible through the view. The following check options are supported:

`LOCAL`

New rows are only checked against the conditions defined directly in the view itself. Any conditions defined on underlying base views are not checked (unless they also specify the `CHECK OPTION`).

`CASCADED`

New rows are checked against the conditions of the view and all underlying base views. If the `CHECK OPTION` is specified, and neither `LOCAL` nor `CASCADED` is specified, then `CASCADED` is assumed.

The `CHECK OPTION` may not be used with `RECURSIVE` views.

Note that the `CHECK OPTION` is only supported on views that are automatically updatable, and do not have `INSTEAD OF` triggers or `INSTEAD` rules. If an automatically updatable view is defined on top of a base view that has `INSTEAD OF` triggers, then the `LOCAL CHECK OPTION` may be used to check the conditions on the automatically updatable view, but the conditions on the base view with `INSTEAD OF` triggers will not be checked (a cascaded check option will not cascade down to a trigger-updatable view, and any check options defined directly on a trigger-updatable view will be ignored). If the view or any of its base relations has an `INSTEAD` rule that causes the `INSERT` or `UPDATE` command to be rewritten, then all check options will be ignored in the rewritten query, including any checks from automatically updatable views defined on top of the relation with the `INSTEAD` rule. `MERGE` is not supported if the view or any of its base relations have rules.

Notes

Use the [DROP VIEW](#) statement to drop views.

Be careful that the names and types of the view's columns will be assigned the way you want. For example:

```
CREATE VIEW vista AS SELECT 'Hello World';
```

is bad form because the column name defaults to `?column?`; also, the column data type defaults to `text`, which might not be what you wanted. Better style for a string literal in a view's result is something like:

```
CREATE VIEW vista AS SELECT text 'Hello World' AS hello;
```

By default, access to the underlying base relations referenced in the view is determined by the permissions of the view owner. In some cases, this can be used to provide secure but restricted access to the underlying tables. However, not all views are secure against tampering; see [Section 42.5](#) for details.

If the view has the `security_invoker` property set to `true`, access to the underlying base relations is determined by the permissions of the user executing the query, rather than the view owner. Thus, the user of a security invoker view must have the relevant permissions on the view and its underlying base relations.

If any of the underlying base relations is a security invoker view, it will be treated as if it had been accessed directly from the original query. Thus, a security invoker view will always check its underlying base relations using the permissions of the current user, even if it is accessed from a view without the `security_invoker` property.

If any of the underlying base relations has [row-level security](#) enabled, then by default, the row-level security policies of the view owner are applied, and access to any additional relations referred to by those policies is determined by the permissions of the view owner. However, if the view has `security_invoker` set to `true`, then the policies and permissions of the invoking user are used instead, as if the base relations had been referenced directly from the query using the view.

Functions called in the view are treated the same as if they had been called directly from the query using the view. Therefore, the user of a view must have permissions to call all functions used by the view. Functions in the view are executed with the privileges of the user executing the query or the function owner, depending on whether the functions are defined as `SECURITY INVOKER` or `SECURITY DEFINER`. Thus, for example, calling `CURRENT_USER` directly in a view will always return the invoking user, not the view owner. This is not affected by the view's `security_invoker` setting, and so a view with `security_invoker` set to `false` is *not* equivalent to a `SECURITY DEFINER` function and those concepts should not be confused.

The user creating or replacing a view must have `USAGE` privileges on any schemas referred to in the view query, in order to look up the referenced objects in those schemas. Note, however, that this lookup only happens when the view is created or replaced. Therefore, the user of the view only requires the `USAGE` privilege on the schema containing the view, not on the schemas referred to in the view query, even for a security invoker view.

When `CREATE OR REPLACE VIEW` is used on an existing view, only the view's defining `SELECT` rule, plus any `WITH (...)` parameters and its `CHECK OPTION` are changed. Other view properties, including ownership, permissions, and non-`SELECT` rules, remain unchanged. You must own the view to replace it (this includes being a member of the owning role).

Updatable Views

Simple views are automatically updatable: the system will allow `INSERT`, `UPDATE`, `DELETE`, and `MERGE` statements to be used on the view in the same way as on a regular table. A view is automatically updatable if it satisfies all of the following conditions:

- The view must have exactly one entry in its `FROM` list, which must be a table or another updatable view.
- The view definition must not contain `WITH`, `DISTINCT`, `GROUP BY`, `HAVING`, `LIMIT`, or `OFFSET` clauses at the top level.
- The view definition must not contain set operations (`UNION`, `INTERSECT` or `EXCEPT`) at the top level.

- The view's select list must not contain any aggregates, window functions or set-returning functions.

An automatically updatable view may contain a mix of updatable and non-updatable columns. A column is updatable if it is a simple reference to an updatable column of the underlying base relation; otherwise the column is read-only, and an error will be raised if an `INSERT`, `UPDATE`, or `MERGE` statement attempts to assign a value to it.

If the view is automatically updatable the system will convert any `INSERT`, `UPDATE`, `DELETE`, or `MERGE` statement on the view into the corresponding statement on the underlying base relation. `INSERT` statements that have an `ON CONFLICT UPDATE` clause are fully supported.

If an automatically updatable view contains a `WHERE` condition, the condition restricts which rows of the base relation are available to be modified by `UPDATE`, `DELETE`, and `MERGE` statements on the view. However, an `UPDATE` or `MERGE` is allowed to change a row so that it no longer satisfies the `WHERE` condition, and thus is no longer visible through the view. Similarly, an `INSERT` or `MERGE` command can potentially insert base-relation rows that do not satisfy the `WHERE` condition and thus are not visible through the view (`ON CONFLICT UPDATE` may similarly affect an existing row not visible through the view). The `CHECK OPTION` may be used to prevent `INSERT`, `UPDATE`, and `MERGE` commands from creating such rows that are not visible through the view.

If an automatically updatable view is marked with the `security_barrier` property then all the view's `WHERE` conditions (and any conditions using operators which are marked as `LEAKPROOF`) will always be evaluated before any conditions that a user of the view has added. See [Section 42.5](#) for full details. Note that, due to this, rows which are not ultimately returned (because they do not pass the user's `WHERE` conditions) may still end up being locked. `EXPLAIN` can be used to see which conditions are applied at the relation level (and therefore do not lock rows) and which are not.

A more complex view that does not satisfy all these conditions is read-only by default: the system will not allow an `INSERT`, `UPDATE`, `DELETE`, or `MERGE` on the view. You can get the effect of an updatable view by creating `INSTEAD OF` triggers on the view, which must convert attempted inserts, etc. on the view into appropriate actions on other tables. For more information see [CREATE TRIGGER](#). Another possibility is to create rules (see [CREATE RULE](#)), but in practice triggers are easier to understand and use correctly. Also note that `MERGE` is not supported on relations with rules.

Note that the user performing the insert, update or delete on the view must have the corresponding insert, update or delete privilege on the view. In addition, by default, the view's owner must have the relevant privileges on the underlying base relations, whereas the user performing the update does not need any permissions on the underlying base relations (see [Section 42.5](#)). However, if the view has `security_invoker` set to `true`, the user performing the update, rather than the view owner, must have the relevant privileges on the underlying base relations.

Examples

Create a view consisting of all comedy films:

```
CREATE VIEW comedies AS
  SELECT *
  FROM films
  WHERE kind = 'Comedy';
```

This will create a view containing the columns that are in the `film` table at the time of view creation. Though `*` was used to create the view, columns added later to the table will not be part of the view.

Create a view with `LOCAL CHECK OPTION`:

```
CREATE VIEW universal_comedies AS
  SELECT *
  FROM comedies
  WHERE classification = 'U'
  WITH LOCAL CHECK OPTION;
```


This will create a view based on the `comedies` view, showing only films with `kind = 'Comedy'` and `classification = 'U'`. Any attempt to `INSERT` or `UPDATE` a row in the view will be rejected if the new row doesn't have `classification = 'U'`, but the film `kind` will not be checked.

Create a view with `CASCADDED CHECK OPTION`:

```
CREATE VIEW pg_comedies AS
  SELECT *
  FROM comedies
  WHERE classification = 'PG'
  WITH CASCADDED CHECK OPTION;
```

This will create a view that checks both the `kind` and `classification` of new rows.

Create a view with a mix of updatable and non-updatable columns:

```
CREATE VIEW comedies AS
  SELECT f.*,
         country_code_to_name(f.country_code) AS country,
         (SELECT avg(r.rating)
          FROM user_ratings r
          WHERE r.film_id = f.id) AS avg_rating
  FROM films f
  WHERE f.kind = 'Comedy';
```

This view will support `INSERT`, `UPDATE` and `DELETE`. All the columns from the `films` table will be updatable, whereas the computed columns `country` and `avg_rating` will be read-only.

Create a recursive view consisting of the numbers from 1 to 100:

```
CREATE RECURSIVE VIEW public.nums_1_100 (n) AS
  VALUES (1)
  UNION ALL
  SELECT n+1 FROM nums_1_100 WHERE n < 100;
```

Notice that although the recursive view's name is schema-qualified in this `CREATE`, its internal self-reference is not schema-qualified. This is because the implicitly-created CTE's name cannot be schema-qualified.

Compatibility

`CREATE OR REPLACE VIEW` is a Postgres Pro language extension. So is the concept of a temporary view. The `WITH (...)` clause is an extension as well, as are security barrier views and security invoker views.

See Also

[ALTER VIEW](#), [DROP VIEW](#), [CREATE MATERIALIZED VIEW](#)

DEALLOCATE

DEALLOCATE — deallocate a prepared statement

Synopsis

```
DEALLOCATE [ PREPARE ] { name | ALL }
```

Description

DEALLOCATE is used to deallocate a previously prepared SQL statement. If you do not explicitly deallocate a prepared statement, it is deallocated when the session ends.

For more information on prepared statements, see [PREPARE](#).

Parameters

PREPARE

This key word is ignored.

name

The name of the prepared statement to deallocate.

ALL

Deallocate all prepared statements.

Compatibility

The SQL standard includes a DEALLOCATE statement, but it is only for use in embedded SQL.

See Also

[EXECUTE](#), [PREPARE](#)

DECLARE

DECLARE — define a cursor

Synopsis

```
DECLARE name [ BINARY ] [ ASENSITIVE | INSENSITIVE ] [ [ NO ] SCROLL ]  
CURSOR [ { WITH | WITHOUT } HOLD ] FOR query
```

Description

DECLARE allows a user to create cursors, which can be used to retrieve a small number of rows at a time out of a larger query. After the cursor is created, rows are fetched from it using [FETCH](#).

Note

This page describes usage of cursors at the SQL command level. If you are trying to use cursors inside a PL/pgSQL function, the rules are different — see [Section 44.7](#).

Parameters

name

The name of the cursor to be created. This must be different from any other active cursor name in the session.

BINARY

Causes the cursor to return data in binary rather than in text format.

ASENSITIVE

INSENSITIVE

Cursor sensitivity determines whether changes to the data underlying the cursor, done in the same transaction, after the cursor has been declared, are visible in the cursor. `INSENSITIVE` means they are not visible, `ASENSITIVE` means the behavior is implementation-dependent. A third behavior, `SENSITIVE`, meaning that such changes are visible in the cursor, is not available in Postgres Pro. In Postgres Pro, all cursors are insensitive; so these key words have no effect and are only accepted for compatibility with the SQL standard.

Specifying `INSENSITIVE` together with `FOR UPDATE` or `FOR SHARE` is an error.

SCROLL

NO SCROLL

`SCROLL` specifies that the cursor can be used to retrieve rows in a nonsequential fashion (e.g., backward). Depending upon the complexity of the query's execution plan, specifying `SCROLL` might impose a performance penalty on the query's execution time. `NO SCROLL` specifies that the cursor cannot be used to retrieve rows in a nonsequential fashion. The default is to allow scrolling in some cases; this is not the same as specifying `SCROLL`. See [Notes](#) below for details.

WITH HOLD

WITHOUT HOLD

`WITH HOLD` specifies that the cursor can continue to be used after the transaction that created it successfully commits. `WITHOUT HOLD` specifies that the cursor cannot be used outside of the transaction that created it. If neither `WITHOUT HOLD` nor `WITH HOLD` is specified, `WITHOUT HOLD` is the default.

query

A `SELECT` or `VALUES` command which will provide the rows to be returned by the cursor.

The key words `ASENSITIVE`, `BINARY`, `INSENSITIVE`, and `SCROLL` can appear in any order.

Notes

Normal cursors return data in text format, the same as a `SELECT` would produce. The `BINARY` option specifies that the cursor should return data in binary format. This reduces conversion effort for both the server and client, at the cost of more programmer effort to deal with platform-dependent binary data formats. As an example, if a query returns a value of one from an integer column, you would get a string of 1 with a default cursor, whereas with a binary cursor you would get a 4-byte field containing the internal representation of the value (in big-endian byte order).

Binary cursors should be used carefully. Many applications, including `psql`, are not prepared to handle binary cursors and expect data to come back in the text format.

Note

When the client application uses the “extended query” protocol to issue a `FETCH` command, the `Bind` protocol message specifies whether data is to be retrieved in text or binary format. This choice overrides the way that the cursor is defined. The concept of a binary cursor as such is thus obsolete when using extended query protocol — any cursor can be treated as either text or binary.

Unless `WITH HOLD` is specified, the cursor created by this command can only be used within the current transaction. Thus, `DECLARE` without `WITH HOLD` is useless outside a transaction block: the cursor would survive only to the completion of the statement. Therefore Postgres Pro reports an error if such a command is used outside a transaction block. Use `BEGIN` and `COMMIT` (or `ROLLBACK`) to define a transaction block.

If `WITH HOLD` is specified and the transaction that created the cursor successfully commits, the cursor can continue to be accessed by subsequent transactions in the same session. (But if the creating transaction is aborted, the cursor is removed.) A cursor created with `WITH HOLD` is closed when an explicit `CLOSE` command is issued on it, or the session ends. In the current implementation, the rows represented by a held cursor are copied into a temporary file or memory area so that they remain available for subsequent transactions.

`WITH HOLD` may not be specified when the query includes `FOR UPDATE` or `FOR SHARE`.

The `SCROLL` option should be specified when defining a cursor that will be used to fetch backwards. This is required by the SQL standard. However, for compatibility with earlier versions, Postgres Pro will allow backward fetches without `SCROLL`, if the cursor's query plan is simple enough that no extra overhead is needed to support it. However, application developers are advised not to rely on using backward fetches from a cursor that has not been created with `SCROLL`. If `NO SCROLL` is specified, then backward fetches are disallowed in any case.

Backward fetches are also disallowed when the query includes `FOR UPDATE` or `FOR SHARE`; therefore `SCROLL` may not be specified in this case.

Caution

Scrollable cursors may give unexpected results if they invoke any volatile functions (see [Section 39.7](#)). When a previously fetched row is re-fetched, the functions might be re-executed, perhaps leading to results different from the first time. It's best to specify `NO SCROLL` for a query involving volatile functions. If that is not practical, one workaround is to declare the cursor `SCROLL WITH HOLD` and commit the transaction before reading any rows from it. This will force the entire

output of the cursor to be materialized in temporary storage, so that volatile functions are executed exactly once for each row.

If the cursor's query includes `FOR UPDATE` or `FOR SHARE`, then returned rows are locked at the time they are first fetched, in the same way as for a regular [SELECT](#) command with these options. In addition, the returned rows will be the most up-to-date versions.

Caution

It is generally recommended to use `FOR UPDATE` if the cursor is intended to be used with `UPDATE ... WHERE CURRENT OF` or `DELETE ... WHERE CURRENT OF`. Using `FOR UPDATE` prevents other sessions from changing the rows between the time they are fetched and the time they are updated. Without `FOR UPDATE`, a subsequent `WHERE CURRENT OF` command will have no effect if the row was changed since the cursor was created.

Another reason to use `FOR UPDATE` is that without it, a subsequent `WHERE CURRENT OF` might fail if the cursor query does not meet the SQL standard's rules for being “simply updatable” (in particular, the cursor must reference just one table and not use grouping or `ORDER BY`). Cursors that are not simply updatable might work, or might not, depending on plan choice details; so in the worst case, an application might work in testing and then fail in production. If `FOR UPDATE` is specified, the cursor is guaranteed to be updatable.

The main reason not to use `FOR UPDATE` with `WHERE CURRENT OF` is if you need the cursor to be scrollable, or to be isolated from concurrent updates (that is, continue to show the old data). If this is a requirement, pay close heed to the caveats shown above.

The SQL standard only makes provisions for cursors in embedded SQL. The Postgres Pro server does not implement an `OPEN` statement for cursors; a cursor is considered to be open when it is declared. However, ECPG, the embedded SQL preprocessor for Postgres Pro, supports the standard SQL cursor conventions, including those involving `DECLARE` and `OPEN` statements.

The server data structure underlying an open cursor is called a *portal*. Portal names are exposed in the client protocol: a client can fetch rows directly from an open portal, if it knows the portal name. When creating a cursor with `DECLARE`, the portal name is the same as the cursor name.

You can see all available cursors by querying the [pg_cursors](#) system view.

Examples

To declare a cursor:

```
DECLARE liahona CURSOR FOR SELECT * FROM films;
```

See [FETCH](#) for more examples of cursor usage.

Compatibility

The SQL standard allows cursors only in embedded SQL and in modules. Postgres Pro permits cursors to be used interactively.

According to the SQL standard, changes made to insensitive cursors by `UPDATE ... WHERE CURRENT OF` and `DELETE ... WHERE CURRENT OF` statements are visible in that same cursor. Postgres Pro treats these statements like all other data changing statements in that they are not visible in insensitive cursors.

Binary cursors are a Postgres Pro extension.

See Also

[CLOSE](#), [FETCH](#), [MOVE](#)

DELETE

DELETE — delete rows of a table

Synopsis

```
[ WITH [ RECURSIVE ] with_query [, ...] ]  
DELETE FROM [ ONLY ] table_name [ * ] [ [ AS ] alias ]  
    [ USING from_item [, ...] ]  
    [ WHERE condition | WHERE CURRENT OF cursor_name ]  
    [ RETURNING { * | output_expression [ [ AS ] output_name ] } [, ...] ]
```

Description

DELETE deletes rows that satisfy the WHERE clause from the specified table. If the WHERE clause is absent, the effect is to delete all rows in the table. The result is a valid, but empty table.

Tip

TRUNCATE provides a faster mechanism to remove all rows from a table.

There are two ways to delete rows in a table using information contained in other tables in the database: using sub-selects, or specifying additional tables in the USING clause. Which technique is more appropriate depends on the specific circumstances.

The optional RETURNING clause causes DELETE to compute and return value(s) based on each row actually deleted. Any expression using the table's columns, and/or columns of other tables mentioned in USING, can be computed. The syntax of the RETURNING list is identical to that of the output list of SELECT.

You must have the DELETE privilege on the table to delete from it, as well as the SELECT privilege for any table in the USING clause or whose values are read in the *condition*.

Parameters

with_query

The WITH clause allows you to specify one or more subqueries that can be referenced by name in the DELETE query. See [Section 7.8](#) and [SELECT](#) for details.

table_name

The name (optionally schema-qualified) of the table to delete rows from. If ONLY is specified before the table name, matching rows are deleted from the named table only. If ONLY is not specified, matching rows are also deleted from any tables inheriting from the named table. Optionally, * can be specified after the table name to explicitly indicate that descendant tables are included.

alias

A substitute name for the target table. When an alias is provided, it completely hides the actual name of the table. For example, given DELETE FROM foo AS f, the remainder of the DELETE statement must refer to this table as f not foo.

from_item

A table expression allowing columns from other tables to appear in the WHERE condition. This uses the same syntax as the FROM clause of a SELECT statement; for example, an alias for the table name can be specified. Do not repeat the target table as a *from_item* unless you wish to set up a self-join (in which case it must appear with an alias in the *from_item*).

condition

An expression that returns a value of type `boolean`. Only rows for which this expression returns `true` will be deleted.

cursor_name

The name of the cursor to use in a `WHERE CURRENT OF` condition. The row to be deleted is the one most recently fetched from this cursor. The cursor must be a non-grouping query on the `DELETE`'s target table. Note that `WHERE CURRENT OF` cannot be specified together with a Boolean condition. See [DECLARE](#) for more information about using cursors with `WHERE CURRENT OF`.

output_expression

An expression to be computed and returned by the `DELETE` command after each row is deleted. The expression can use any column names of the table named by *table_name* or table(s) listed in `USING`. Write `*` to return all columns.

output_name

A name to use for a returned column.

Outputs

On successful completion, a `DELETE` command returns a command tag of the form

```
DELETE count
```

The *count* is the number of rows deleted. Note that the number may be less than the number of rows that matched the *condition* when deletes were suppressed by a `BEFORE DELETE` trigger. If *count* is 0, no rows were deleted by the query (this is not considered an error).

If the `DELETE` command contains a `RETURNING` clause, the result will be similar to that of a `SELECT` statement containing the columns and values defined in the `RETURNING` list, computed over the row(s) deleted by the command.

Notes

Postgres Pro lets you reference columns of other tables in the `WHERE` condition by specifying the other tables in the `USING` clause. For example, to delete all films produced by a given producer, one can do:

```
DELETE FROM films USING producers
WHERE producer_id = producers.id AND producers.name = 'foo';
```

What is essentially happening here is a join between `films` and `producers`, with all successfully joined `films` rows being marked for deletion. This syntax is not standard. A more standard way to do it is:

```
DELETE FROM films
WHERE producer_id IN (SELECT id FROM producers WHERE name = 'foo');
```

In some cases the join style is easier to write or faster to execute than the sub-select style.

Examples

Delete all films but musicals:

```
DELETE FROM films WHERE kind <> 'Musical';
```

Clear the table `films`:

```
DELETE FROM films;
```

Delete completed tasks, returning full details of the deleted rows:

```
DELETE FROM tasks WHERE status = 'DONE' RETURNING *;
```

Delete the row of `tasks` on which the cursor `c_tasks` is currently positioned:

```
DELETE FROM tasks WHERE CURRENT OF c_tasks;
```

While there is no `LIMIT` clause for `DELETE`, it is possible to get a similar effect using the same method described in [the documentation of `UPDATE`](#):

```
WITH delete_batch AS (  
  SELECT l.ctid FROM user_logs AS l  
    WHERE l.status = 'archived'  
    ORDER BY l.creation_date  
  FOR UPDATE  
  LIMIT 10000  
)  
DELETE FROM user_logs AS dl  
  USING delete_batch AS del  
  WHERE dl.ctid = del.ctid;
```

Compatibility

This command conforms to the SQL standard, except that the `USING` and `RETURNING` clauses are Postgres Pro extensions, as is the ability to use `WITH` with `DELETE`.

See Also

[TRUNCATE](#)

DISCARD

DISCARD — discard session state

Synopsis

```
DISCARD { ALL | PLANS | SEQUENCES | TEMPORARY | TEMP }
```

Description

`DISCARD` releases internal resources associated with a database session. This command is useful for partially or fully resetting the session's state. There are several subcommands to release different types of resources; the `DISCARD ALL` variant subsumes all the others, and also resets additional state.

Parameters

`PLANS`

Releases all cached query plans, forcing re-planning to occur the next time the associated prepared statement is used.

`SEQUENCES`

Discards all cached sequence-related state, including `currval()`/`lastval()` information and any pre-allocated sequence values that have not yet been returned by `nextval()`. (See [CREATE SEQUENCE](#) for a description of preallocated sequence values.)

`TEMPORARY` or `TEMP`

Drops all temporary tables created in the current session.

`ALL`

Releases all temporary resources associated with the current session and resets the session to its initial state. Currently, this has the same effect as executing the following sequence of statements:

```
CLOSE ALL;  
SET SESSION AUTHORIZATION DEFAULT;  
RESET ALL;  
DEALLOCATE ALL;  
UNLISTEN *;  
SELECT pg_advisory_unlock_all();  
DISCARD PLANS;  
DISCARD TEMP;  
DISCARD SEQUENCES;
```

Notes

`DISCARD ALL` cannot be executed inside a transaction block.

Compatibility

`DISCARD` is a Postgres Pro extension.

DO

DO — execute an anonymous code block

Synopsis

```
DO [ LANGUAGE lang_name ] code
```

Description

DO executes an anonymous code block, or in other words a transient anonymous function in a procedural language.

The code block is treated as though it were the body of a function with no parameters, returning `void`. It is parsed and executed a single time.

The optional `LANGUAGE` clause can be written either before or after the code block.

Parameters

code

The procedural language code to be executed. This must be specified as a string literal, just as in `CREATE FUNCTION`. Use of a dollar-quoted literal is recommended.

lang_name

The name of the procedural language the code is written in. If omitted, the default is `plpgsql`.

Notes

The procedural language to be used must already have been installed into the current database by means of `CREATE EXTENSION`. `plpgsql` is installed by default, but other languages are not.

The user must have `USAGE` privilege for the procedural language, or must be a superuser if the language is untrusted. This is the same privilege requirement as for creating a function in the language.

If `DO` is executed in a transaction block, then the procedure code cannot execute transaction control statements. Transaction control statements are only allowed if `DO` is executed in its own transaction.

Examples

Grant all privileges on all views in schema `public` to role `webuser`:

```
DO $$DECLARE r record;
BEGIN
    FOR r IN SELECT table_schema, table_name FROM information_schema.tables
              WHERE table_type = 'VIEW' AND table_schema = 'public'
    LOOP
        EXECUTE 'GRANT ALL ON ' || quote_ident(r.table_schema) || '.' ||
quote_ident(r.table_name) || ' TO webuser';
    END LOOP;
END$$;
```

Compatibility

There is no `DO` statement in the SQL standard.

See Also

[CREATE LANGUAGE](#)

DROP ACCESS METHOD

DROP ACCESS METHOD — remove an access method

Synopsis

```
DROP ACCESS METHOD [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP ACCESS METHOD removes an existing access method. Only superusers can drop access methods.

Parameters

IF EXISTS

Do not throw an error if the access method does not exist. A notice is issued in this case.

name

The name of an existing access method.

CASCADE

Automatically drop objects that depend on the access method (such as operator classes, operator families, and indexes), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the access method if any objects depend on it. This is the default.

Examples

Drop the access method `heptree`:

```
DROP ACCESS METHOD heptree;
```

Compatibility

DROP ACCESS METHOD is a Postgres Pro extension.

See Also

[CREATE ACCESS METHOD](#)

DROP AGGREGATE

DROP AGGREGATE — remove an aggregate function

Synopsis

```
DROP AGGREGATE [ IF EXISTS ] name ( aggregate_signature ) [, ...] [ CASCADE |  
RESTRICT ]
```

where *aggregate_signature* is:

```
* |  
[ argmode ] [ argname ] argtype [ , ... ] |  
[ [ argmode ] [ argname ] argtype [ , ... ] ] ORDER BY [ argmode ] [ argname ] argtype  
[ , ... ]
```

Description

DROP AGGREGATE removes an existing aggregate function. To execute this command the current user must be the owner of the aggregate function.

Parameters

IF EXISTS

Do not throw an error if the aggregate does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing aggregate function.

argmode

The mode of an argument: IN or VARIADIC. If omitted, the default is IN.

argname

The name of an argument. Note that DROP AGGREGATE does not actually pay any attention to argument names, since only the argument data types are needed to determine the aggregate function's identity.

argtype

An input data type on which the aggregate function operates. To reference a zero-argument aggregate function, write * in place of the list of argument specifications. To reference an ordered-set aggregate function, write ORDER BY between the direct and aggregated argument specifications.

CASCADE

Automatically drop objects that depend on the aggregate function (such as views using it), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the aggregate function if any objects depend on it. This is the default.

Notes

Alternative syntaxes for referencing ordered-set aggregates are described under [ALTER AGGREGATE](#).

Examples

To remove the aggregate function `myavg` for type `integer`:

```
DROP AGGREGATE myavg(integer);
```

To remove the hypothetical-set aggregate function `myrank`, which takes an arbitrary list of ordering columns and a matching list of direct arguments:

```
DROP AGGREGATE myrank(VARIADIC "any" ORDER BY VARIADIC "any");
```

To remove multiple aggregate functions in one command:

```
DROP AGGREGATE myavg(integer), myavg(bigint);
```

Compatibility

There is no `DROP AGGREGATE` statement in the SQL standard.

See Also

[ALTER AGGREGATE](#), [CREATE AGGREGATE](#)

DROP CAST

DROP CAST — remove a cast

Synopsis

```
DROP CAST [ IF EXISTS ] (source_type AS target_type) [ CASCADE | RESTRICT ]
```

Description

DROP CAST removes a previously defined cast.

To be able to drop a cast, you must own the source or the target data type. These are the same privileges that are required to create a cast.

Parameters

IF EXISTS

Do not throw an error if the cast does not exist. A notice is issued in this case.

source_type

The name of the source data type of the cast.

target_type

The name of the target data type of the cast.

CASCADE

RESTRICT

These key words do not have any effect, since there are no dependencies on casts.

Examples

To drop the cast from type `text` to type `int`:

```
DROP CAST (text AS int);
```

Compatibility

The DROP CAST command conforms to the SQL standard.

See Also

[CREATE CAST](#)

DROP COLLATION

DROP COLLATION — remove a collation

Synopsis

```
DROP COLLATION [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP COLLATION removes a previously defined collation. To be able to drop a collation, you must own the collation.

Parameters

IF EXISTS

Do not throw an error if the collation does not exist. A notice is issued in this case.

name

The name of the collation. The collation name can be schema-qualified.

CASCADE

Automatically drop objects that depend on the collation, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the collation if any objects depend on it. This is the default.

Examples

To drop the collation named `german`:

```
DROP COLLATION german;
```

Compatibility

The DROP COLLATION command conforms to the SQL standard, apart from the IF EXISTS option, which is a Postgres Pro extension.

See Also

[ALTER COLLATION](#), [CREATE COLLATION](#)

DROP CONVERSION

DROP CONVERSION — remove a conversion

Synopsis

```
DROP CONVERSION [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

`DROP CONVERSION` removes a previously defined conversion. To be able to drop a conversion, you must own the conversion.

Parameters

`IF EXISTS`

Do not throw an error if the conversion does not exist. A notice is issued in this case.

name

The name of the conversion. The conversion name can be schema-qualified.

`CASCADE`

`RESTRICT`

These key words do not have any effect, since there are no dependencies on conversions.

Examples

To drop the conversion named `myname`:

```
DROP CONVERSION myname;
```

Compatibility

There is no `DROP CONVERSION` statement in the SQL standard, but a `DROP TRANSLATION` statement that goes along with the `CREATE TRANSLATION` statement that is similar to the `CREATE CONVERSION` statement in Postgres Pro.

See Also

[ALTER CONVERSION](#), [CREATE CONVERSION](#)

DROP DATABASE

DROP DATABASE — remove a database

Synopsis

```
DROP DATABASE [ IF EXISTS ] name [ [ WITH ] ( option [, ...] ) ]
```

where *option* can be:

FORCE

Description

DROP DATABASE drops a database. It removes the catalog entries for the database and deletes the directory containing the data. It can only be executed by the database owner. It cannot be executed while you are connected to the target database. (Connect to `postgres` or any other database to issue this command.) Also, if anyone else is connected to the target database, this command will fail unless you use the `FORCE` option described below.

DROP DATABASE cannot be undone. Use it with care!

Parameters

IF EXISTS

Do not throw an error if the database does not exist. A notice is issued in this case.

name

The name of the database to remove.

FORCE

Attempt to terminate all existing connections to the target database. It doesn't terminate if prepared transactions, active logical replication slots or subscriptions are present in the target database.

This terminates background worker connections and connections that the current user has permission to terminate with `pg_terminate_backend`, described in [Section 9.28.2](#). If connections would remain, this command will fail.

Notes

DROP DATABASE cannot be executed inside a transaction block.

This command cannot be executed while connected to the target database. Thus, it might be more convenient to use the program [dropdb](#) instead, which is a wrapper around this command.

Compatibility

There is no DROP DATABASE statement in the SQL standard.

See Also

[CREATE DATABASE](#)

DROP DOMAIN

DROP DOMAIN — remove a domain

Synopsis

```
DROP DOMAIN [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP DOMAIN removes a domain. Only the owner of a domain can remove it.

Parameters

IF EXISTS

Do not throw an error if the domain does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing domain.

CASCADE

Automatically drop objects that depend on the domain (such as table columns), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the domain if any objects depend on it. This is the default.

Examples

To remove the domain `box`:

```
DROP DOMAIN box;
```

Compatibility

This command conforms to the SQL standard, except for the IF EXISTS option, which is a Postgres Pro extension.

See Also

[CREATE DOMAIN](#), [ALTER DOMAIN](#)

DROP EVENT TRIGGER

DROP EVENT TRIGGER — remove an event trigger

Synopsis

```
DROP EVENT TRIGGER [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP EVENT TRIGGER removes an existing event trigger. To execute this command, the current user must be the owner of the event trigger.

Parameters

IF EXISTS

Do not throw an error if the event trigger does not exist. A notice is issued in this case.

name

The name of the event trigger to remove.

CASCADE

Automatically drop objects that depend on the trigger, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the trigger if any objects depend on it. This is the default.

Examples

Destroy the trigger `snitch`:

```
DROP EVENT TRIGGER snitch;
```

Compatibility

There is no DROP EVENT TRIGGER statement in the SQL standard.

See Also

[CREATE EVENT TRIGGER](#), [ALTER EVENT TRIGGER](#)

DROP EXTENSION

DROP EXTENSION — remove an extension

Synopsis

```
DROP EXTENSION [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP EXTENSION removes extensions from the database. Dropping an extension causes its member objects, and other explicitly dependent routines (see [ALTER ROUTINE](#), the `DEPENDS ON EXTENSION extension_name` action), to be dropped as well.

You must own the extension to use DROP EXTENSION.

Parameters

IF EXISTS

Do not throw an error if the extension does not exist. A notice is issued in this case.

name

The name of an installed extension.

CASCADE

Automatically drop objects that depend on the extension, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

This option prevents the specified extensions from being dropped if other objects, besides these extensions, their members, and their explicitly dependent routines, depend on them. This is the default.

Examples

To remove the extension `hstore` from the current database:

```
DROP EXTENSION hstore;
```

This command will fail if any of `hstore`'s objects are in use in the database, for example if any tables have columns of the `hstore` type. Add the `CASCADE` option to forcibly remove those dependent objects as well.

Compatibility

DROP EXTENSION is a Postgres Pro extension.

See Also

[CREATE EXTENSION](#), [ALTER EXTENSION](#)

DROP FUNCTION

DROP FUNCTION — remove a function

Synopsis

```
DROP FUNCTION [ IF EXISTS ] name [ ( [ argmode ] [ argname ] argtype [, ...] ) ]  
[ , ... ]  
[ CASCADE | RESTRICT ]
```

Description

DROP FUNCTION removes the definition of an existing function. To execute this command the user must be the owner of the function. The argument types to the function must be specified, since several different functions can exist with the same name and different argument lists.

Parameters

IF EXISTS

Do not throw an error if the function does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing function. If no argument list is specified, the name must be unique in its schema.

argmode

The mode of an argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN. Note that DROP FUNCTION does not actually pay any attention to OUT arguments, since only the input arguments are needed to determine the function's identity. So it is sufficient to list the IN, INOUT, and VARIADIC arguments.

argname

The name of an argument. Note that DROP FUNCTION does not actually pay any attention to argument names, since only the argument data types are needed to determine the function's identity.

argtype

The data type(s) of the function's arguments (optionally schema-qualified), if any.

CASCADE

Automatically drop objects that depend on the function (such as operators or triggers), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the function if any objects depend on it. This is the default.

Examples

This command removes the square root function:

```
DROP FUNCTION sqrt(integer);
```

Drop multiple functions in one command:

```
DROP FUNCTION sqrt(integer), sqrt(bigint);
```

If the function name is unique in its schema, it can be referred to without an argument list:

```
DROP FUNCTION update_employee_salaries;
```

Note that this is different from

```
DROP FUNCTION update_employee_salaries();
```

which refers to a function with zero arguments, whereas the first variant can refer to a function with any number of arguments, including zero, as long as the name is unique.

Compatibility

This command conforms to the SQL standard, with these Postgres Pro extensions:

- The standard only allows one function to be dropped per command.
- The `IF EXISTS` option
- The ability to specify argument modes and names

See Also

[CREATE FUNCTION](#), [ALTER FUNCTION](#), [DROP PROCEDURE](#), [DROP ROUTINE](#)

DROP GROUP

DROP GROUP — remove a database role

Synopsis

```
DROP GROUP [ IF EXISTS ] name [, ...]
```

Description

DROP GROUP is now an alias for [DROP ROLE](#).

Compatibility

There is no DROP GROUP statement in the SQL standard.

See Also

[DROP ROLE](#)

DROP INDEX

DROP INDEX — remove an index

Synopsis

```
DROP INDEX [ CONCURRENTLY ] [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP INDEX drops an existing index from the database system. To execute this command you must be the owner of the index.

Parameters

CONCURRENTLY

Drop the index without locking out concurrent selects, inserts, updates, and deletes on the index's table. A normal DROP INDEX acquires an ACCESS EXCLUSIVE lock on the table, blocking other accesses until the index drop can be completed. With this option, the command instead waits until conflicting transactions have completed.

There are several caveats to be aware of when using this option. Only one index name can be specified, and the CASCADE option is not supported. (Thus, an index that supports a UNIQUE or PRIMARY KEY constraint cannot be dropped this way.) Also, regular DROP INDEX commands can be performed within a transaction block, but DROP INDEX CONCURRENTLY cannot. Lastly, indexes on partitioned tables cannot be dropped using this option.

For temporary tables, DROP INDEX is always non-concurrent, as no other session can access them, and non-concurrent index drop is cheaper.

IF EXISTS

Do not throw an error if the index does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an index to remove.

CASCADE

Automatically drop objects that depend on the index, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the index if any objects depend on it. This is the default.

Examples

This command will remove the index `title_idx`:

```
DROP INDEX title_idx;
```

Compatibility

DROP INDEX is a Postgres Pro language extension. There are no provisions for indexes in the SQL standard.

See Also

[CREATE INDEX](#)

DROP LANGUAGE

DROP LANGUAGE — remove a procedural language

Synopsis

```
DROP [ PROCEDURAL ] LANGUAGE [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP LANGUAGE removes the definition of a previously registered procedural language. You must be a superuser or the owner of the language to use DROP LANGUAGE.

Note

As of PostgreSQL 9.1, most procedural languages have been made into “extensions”, and should therefore be removed with [DROP EXTENSION](#) not DROP LANGUAGE.

Parameters

IF EXISTS

Do not throw an error if the language does not exist. A notice is issued in this case.

name

The name of an existing procedural language.

CASCADE

Automatically drop objects that depend on the language (such as functions in the language), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the language if any objects depend on it. This is the default.

Examples

This command removes the procedural language `plsample`:

```
DROP LANGUAGE plsample;
```

Compatibility

There is no DROP LANGUAGE statement in the SQL standard.

See Also

[ALTER LANGUAGE](#), [CREATE LANGUAGE](#)

DROP MATERIALIZED VIEW

DROP MATERIALIZED VIEW — remove a materialized view

Synopsis

```
DROP MATERIALIZED VIEW [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP MATERIALIZED VIEW drops an existing materialized view. To execute this command you must be the owner of the materialized view.

Parameters

IF EXISTS

Do not throw an error if the materialized view does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of the materialized view to remove.

CASCADE

Automatically drop objects that depend on the materialized view (such as other materialized views, or regular views), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the materialized view if any objects depend on it. This is the default.

Examples

This command will remove the materialized view called `order_summary`:

```
DROP MATERIALIZED VIEW order_summary;
```

Compatibility

DROP MATERIALIZED VIEW is a Postgres Pro extension.

See Also

[CREATE MATERIALIZED VIEW](#), [ALTER MATERIALIZED VIEW](#), [REFRESH MATERIALIZED VIEW](#)

DROP OPERATOR

DROP OPERATOR — remove an operator

Synopsis

```
DROP OPERATOR [ IF EXISTS ] name ( { left_type | NONE } , right_type ) [, ...]  
[ CASCADE | RESTRICT ]
```

Description

DROP OPERATOR drops an existing operator from the database system. To execute this command you must be the owner of the operator.

Parameters

IF EXISTS

Do not throw an error if the operator does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing operator.

left_type

The data type of the operator's left operand; write NONE if the operator has no left operand.

right_type

The data type of the operator's right operand.

CASCADE

Automatically drop objects that depend on the operator (such as views using it), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the operator if any objects depend on it. This is the default.

Examples

Remove the power operator a^b for type integer:

```
DROP OPERATOR ^ (integer, integer);
```

Remove the bitwise-complement prefix operator $\sim b$ for type bit:

```
DROP OPERATOR ~ (none, bit);
```

Remove multiple operators in one command:

```
DROP OPERATOR ~ (none, bit), ^ (integer, integer);
```

Compatibility

There is no DROP OPERATOR statement in the SQL standard.

See Also

[CREATE OPERATOR](#), [ALTER OPERATOR](#)

DROP OPERATOR CLASS

DROP OPERATOR CLASS — remove an operator class

Synopsis

```
DROP OPERATOR CLASS [ IF EXISTS ] name USING index_method [ CASCADE | RESTRICT ]
```

Description

DROP OPERATOR CLASS drops an existing operator class. To execute this command you must be the owner of the operator class.

DROP OPERATOR CLASS does not drop any of the operators or functions referenced by the class. If there are any indexes depending on the operator class, you will need to specify CASCADE for the drop to complete.

Parameters

IF EXISTS

Do not throw an error if the operator class does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing operator class.

index_method

The name of the index access method the operator class is for.

CASCADE

Automatically drop objects that depend on the operator class (such as indexes), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the operator class if any objects depend on it. This is the default.

Notes

DROP OPERATOR CLASS will not drop the operator family containing the class, even if there is nothing else left in the family (in particular, in the case where the family was implicitly created by CREATE OPERATOR CLASS). An empty operator family is harmless, but for the sake of tidiness you might wish to remove the family with DROP OPERATOR FAMILY; or perhaps better, use DROP OPERATOR FAMILY in the first place.

Examples

Remove the B-tree operator class `widget_ops`:

```
DROP OPERATOR CLASS widget_ops USING btree;
```

This command will not succeed if there are any existing indexes that use the operator class. Add CASCADE to drop such indexes along with the operator class.

Compatibility

There is no DROP OPERATOR CLASS statement in the SQL standard.

See Also

[ALTER OPERATOR CLASS](#), [CREATE OPERATOR CLASS](#), [DROP OPERATOR FAMILY](#)

DROP OPERATOR FAMILY

DROP OPERATOR FAMILY — remove an operator family

Synopsis

```
DROP OPERATOR FAMILY [ IF EXISTS ] name USING index_method [ CASCADE | RESTRICT ]
```

Description

DROP OPERATOR FAMILY drops an existing operator family. To execute this command you must be the owner of the operator family.

DROP OPERATOR FAMILY includes dropping any operator classes contained in the family, but it does not drop any of the operators or functions referenced by the family. If there are any indexes depending on operator classes within the family, you will need to specify CASCADE for the drop to complete.

Parameters

IF EXISTS

Do not throw an error if the operator family does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing operator family.

index_method

The name of the index access method the operator family is for.

CASCADE

Automatically drop objects that depend on the operator family, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the operator family if any objects depend on it. This is the default.

Examples

Remove the B-tree operator family `float_ops`:

```
DROP OPERATOR FAMILY float_ops USING btree;
```

This command will not succeed if there are any existing indexes that use operator classes within the family. Add CASCADE to drop such indexes along with the operator family.

Compatibility

There is no DROP OPERATOR FAMILY statement in the SQL standard.

See Also

[ALTER OPERATOR FAMILY](#), [CREATE OPERATOR FAMILY](#), [ALTER OPERATOR CLASS](#), [CREATE OPERATOR CLASS](#), [DROP OPERATOR CLASS](#)

DROP OWNED

DROP OWNED — remove database objects owned by a database role

Synopsis

```
DROP OWNED BY { name | CURRENT_ROLE | CURRENT_USER | SESSION_USER } [, ...] [ CASCADE |  
RESTRICT ]
```

Description

DROP OWNED drops all the objects within the current database that are owned by one of the specified roles. Any privileges granted to the given roles on objects in the current database or on shared objects (databases, tablespaces, configuration parameters) will also be revoked.

Parameters

name

The name of a role whose objects will be dropped, and whose privileges will be revoked.

CASCADE

Automatically drop objects that depend on the affected objects, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the objects owned by a role if any other database objects depend on one of the affected objects. This is the default.

Notes

DROP OWNED is often used to prepare for the removal of one or more roles. Because DROP OWNED only affects the objects in the current database, it is usually necessary to execute this command in each database that contains objects owned by a role that is to be removed.

Using the CASCADE option might make the command recurse to objects owned by other users.

The REASSIGN OWNED command is an alternative that reassigns the ownership of all the database objects owned by one or more roles. However, REASSIGN OWNED does not deal with privileges for other objects.

Databases and tablespaces owned by the role(s) will not be removed.

See [Section 21.4](#) for more discussion.

Compatibility

The DROP OWNED command is a Postgres Pro extension.

See Also

[REASSIGN OWNED](#), [DROP ROLE](#)

DROP PACKAGE

DROP PACKAGE — remove a package

Synopsis

```
DROP PACKAGE [ IF EXISTS ] package_name [, ...] [ CASCADE ]
```

Description

DROP PACKAGE removes packages from the database.

A package can only be dropped by its owner or a superuser. Note that the owner can drop the package (and thereby all contained objects) even if they do not own some of the objects within the package.

Parameters

IF EXISTS

Do not throw an error if the package does not exist. A notice is issued in this case.

package_name

The name of a package.

CASCADE

Automatically drop objects (tables, functions, etc.) that are contained in the package, and in turn all objects that depend on those objects (see [Section 5.15](#)).

Notes

Using the CASCADE option might make the command find objects to be removed in other packages besides the one(s) named. In this case, other packages are dropped as well.

Examples

To remove package `counter` from the database, along with everything it contains and all the dependent objects:

```
DROP PACKAGE counter CASCADE;
```

See Also

[CREATE PACKAGE](#)

DROP POLICY

DROP POLICY — remove a row-level security policy from a table

Synopsis

```
DROP POLICY [ IF EXISTS ] name ON table_name [ CASCADE | RESTRICT ]
```

Description

DROP POLICY removes the specified policy from the table. Note that if the last policy is removed for a table and the table still has row-level security enabled via ALTER TABLE, then the default-deny policy will be used. ALTER TABLE ... DISABLE ROW LEVEL SECURITY can be used to disable row-level security for a table, whether policies for the table exist or not.

Parameters

IF EXISTS

Do not throw an error if the policy does not exist. A notice is issued in this case.

name

The name of the policy to drop.

table_name

The name (optionally schema-qualified) of the table that the policy is on.

CASCADE

RESTRICT

These key words do not have any effect, since there are no dependencies on policies.

Examples

To drop the policy called p1 on the table named my_table:

```
DROP POLICY p1 ON my_table;
```

Compatibility

DROP POLICY is a Postgres Pro extension.

See Also

[CREATE POLICY](#), [ALTER POLICY](#)

DROP PROCEDURE

DROP PROCEDURE — remove a procedure

Synopsis

```
DROP PROCEDURE [ IF EXISTS ] name [ ( [ argmode ] [ argname ] argtype [, ...] ) ]  
[ , ... ]  
[ CASCADE | RESTRICT ]
```

Description

`DROP PROCEDURE` removes the definition of one or more existing procedures. To execute this command the user must be the owner of the procedure(s). The argument types to the procedure(s) usually must be specified, since several different procedures can exist with the same name and different argument lists.

Parameters

`IF EXISTS`

Do not throw an error if the procedure does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing procedure.

argmode

The mode of an argument: `IN`, `OUT`, `INOUT`, or `VARIADIC`. If omitted, the default is `IN` (but see below).

argname

The name of an argument. Note that `DROP PROCEDURE` does not actually pay any attention to argument names, since only the argument data types are used to determine the procedure's identity.

argtype

The data type(s) of the procedure's arguments (optionally schema-qualified), if any. See below for details.

`CASCADE`

Automatically drop objects that depend on the procedure, and in turn all objects that depend on those objects (see [Section 5.15](#)).

`RESTRICT`

Refuse to drop the procedure if any objects depend on it. This is the default.

Notes

If there is only one procedure of the given name, the argument list can be omitted. Omit the parentheses too in this case.

In Postgres Pro, it's sufficient to list the input (including `INOUT`) arguments, because no two routines of the same name are allowed to share the same input-argument list. Moreover, the `DROP` command will not actually check that you wrote the types of `OUT` arguments correctly; so any arguments that are explicitly marked `OUT` are just noise. But writing them is recommendable for consistency with the corresponding `CREATE` command.

For compatibility with the SQL standard, it is also allowed to write all the argument data types (including those of `OUT` arguments) without any *argmode* markers. When this is done, the types of the procedure's

OUT argument(s) *will* be verified against the command. This provision creates an ambiguity, in that when the argument list contains no *argmode* markers, it's unclear which rule is intended. The `DROP` command will attempt the lookup both ways, and will throw an error if two different procedures are found. To avoid the risk of such ambiguity, it's recommendable to write `IN` markers explicitly rather than letting them be defaulted, thus forcing the traditional Postgres Pro interpretation to be used.

The lookup rules just explained are also used by other commands that act on existing procedures, such as `ALTER PROCEDURE` and `COMMENT ON PROCEDURE`.

Examples

If there is only one procedure `do_db_maintenance`, this command is sufficient to drop it:

```
DROP PROCEDURE do_db_maintenance;
```

Given this procedure definition:

```
CREATE PROCEDURE do_db_maintenance(IN target_schema text, OUT results text) ...
```

any one of these commands would work to drop it:

```
DROP PROCEDURE do_db_maintenance(IN target_schema text, OUT results text);
DROP PROCEDURE do_db_maintenance(IN text, OUT text);
DROP PROCEDURE do_db_maintenance(IN text);
DROP PROCEDURE do_db_maintenance(text);
DROP PROCEDURE do_db_maintenance(text, text); -- potentially ambiguous
```

However, the last example would be ambiguous if there is also, say,

```
CREATE PROCEDURE do_db_maintenance(IN target_schema text, IN options text) ...
```

Compatibility

This command conforms to the SQL standard, with these Postgres Pro extensions:

- The standard only allows one procedure to be dropped per command.
- The `IF EXISTS` option is an extension.
- The ability to specify argument modes and names is an extension, and the lookup rules differ when modes are given.

See Also

[CREATE PROCEDURE](#), [ALTER PROCEDURE](#), [DROP FUNCTION](#), [DROP ROUTINE](#)

DROP PROFILE

DROP PROFILE — remove an authentication profile

Synopsis

```
DROP PROFILE [ IF EXISTS ] name
```

Description

The `DROP PROFILE` command removes the specified profile. You must be a database superuser or have the privileges of the `pg_manage_profiles` role to use this command.

A profile cannot be removed if it is still referenced in any database of the cluster; if this is the case, an error is raised. Before dropping the profile, you must drop all the roles to which it is assigned, or assign another profile to them. Use the [ALTER ROLE](#) command if you need to reassign the profile.

Parameters

`IF EXISTS`

Do not throw an error if the profile does not exist. A notice is issued in this case.

name

The name of the profile to remove.

Examples

Drop the specified profile:

```
DROP PROFILE admin_profile;
```

See Also

[CREATE PROFILE](#), [ALTER PROFILE](#), [ALTER ROLE](#)

DROP PUBLICATION

DROP PUBLICATION — remove a publication

Synopsis

```
DROP PUBLICATION [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP PUBLICATION removes an existing publication from the database.

A publication can only be dropped by its owner or a superuser.

Parameters

IF EXISTS

Do not throw an error if the publication does not exist. A notice is issued in this case.

name

The name of an existing publication.

CASCADE

RESTRICT

These key words do not have any effect, since there are no dependencies on publications.

Examples

Drop a publication:

```
DROP PUBLICATION mypublication;
```

Compatibility

DROP PUBLICATION is a Postgres Pro extension.

See Also

[CREATE PUBLICATION](#), [ALTER PUBLICATION](#)

DROP ROLE

DROP ROLE — remove a database role

Synopsis

```
DROP ROLE [ IF EXISTS ] name [, ...]
```

Description

`DROP ROLE` removes the specified role(s). To drop a superuser role, you must be a superuser yourself; to drop non-superuser roles, you must have `CREATEROLE` privilege and have been granted `ADMIN OPTION` on the role.

A role cannot be removed if it is still referenced in any database of the cluster; an error will be raised if so. Before dropping the role, you must drop all the objects it owns (or reassign their ownership) and revoke any privileges the role has been granted on other objects. The [REASSIGN OWNED](#) and [DROP OWNED](#) commands can be useful for this purpose; see [Section 21.4](#) for more discussion.

However, it is not necessary to remove role memberships involving the role; `DROP ROLE` automatically revokes any memberships of the target role in other roles, and of other roles in the target role. The other roles are not dropped nor otherwise affected.

Parameters

`IF EXISTS`

Do not throw an error if the role does not exist. A notice is issued in this case.

name

The name of the role to remove.

Notes

Postgres Pro includes a program [dropuser](#) that has the same functionality as this command (in fact, it calls this command) but can be run from the command shell.

Examples

To drop a role:

```
DROP ROLE jonathan;
```

Compatibility

The SQL standard defines `DROP ROLE`, but it allows only one role to be dropped at a time, and it specifies different privilege requirements than Postgres Pro uses.

See Also

[CREATE ROLE](#), [ALTER ROLE](#), [SET ROLE](#)

DROP ROUTINE

DROP ROUTINE — remove a routine

Synopsis

```
DROP ROUTINE [ IF EXISTS ] name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ]  
[ , ... ]  
[ CASCADE | RESTRICT ]
```

Description

DROP ROUTINE removes the definition of one or more existing routines. The term “routine” includes aggregate functions, normal functions, and procedures. See under [DROP AGGREGATE](#), [DROP FUNCTION](#), and [DROP PROCEDURE](#) for the description of the parameters, more examples, and further details.

Notes

The lookup rules used by DROP ROUTINE are fundamentally the same as for DROP PROCEDURE; in particular, DROP ROUTINE shares that command's behavior of considering an argument list that has no *argmode* markers to be possibly using the SQL standard's definition that OUT arguments are included in the list. (DROP AGGREGATE and DROP FUNCTION do not do that.)

In some cases where the same name is shared by routines of different kinds, it is possible for DROP ROUTINE to fail with an ambiguity error when a more specific command (DROP FUNCTION, etc.) would work. Specifying the argument type list more carefully will also resolve such problems.

These lookup rules are also used by other commands that act on existing routines, such as ALTER ROUTINE and COMMENT ON ROUTINE.

Examples

To drop the routine `foo` for type `integer`:

```
DROP ROUTINE foo(integer);
```

This command will work independent of whether `foo` is an aggregate, function, or procedure.

Compatibility

This command conforms to the SQL standard, with these Postgres Pro extensions:

- The standard only allows one routine to be dropped per command.
- The `IF EXISTS` option is an extension.
- The ability to specify argument modes and names is an extension, and the lookup rules differ when modes are given.
- User-definable aggregate functions are an extension.

See Also

[DROP AGGREGATE](#), [DROP FUNCTION](#), [DROP PROCEDURE](#), [ALTER ROUTINE](#)

Note that there is no `CREATE ROUTINE` command.

DROP RULE

DROP RULE — remove a rewrite rule

Synopsis

```
DROP RULE [ IF EXISTS ] name ON table_name [ CASCADE | RESTRICT ]
```

Description

DROP RULE drops a rewrite rule.

Parameters

IF EXISTS

Do not throw an error if the rule does not exist. A notice is issued in this case.

name

The name of the rule to drop.

table_name

The name (optionally schema-qualified) of the table or view that the rule applies to.

CASCADE

Automatically drop objects that depend on the rule, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the rule if any objects depend on it. This is the default.

Examples

To drop the rewrite rule `newrule`:

```
DROP RULE newrule ON mytable;
```

Compatibility

DROP RULE is a Postgres Pro language extension, as is the entire query rewrite system.

See Also

[CREATE RULE](#), [ALTER RULE](#)

DROP SCHEMA

DROP SCHEMA — remove a schema

Synopsis

```
DROP SCHEMA [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP SCHEMA removes schemas from the database.

A schema can only be dropped by its owner or a superuser. Note that the owner can drop the schema (and thereby all contained objects) even if they do not own some of the objects within the schema.

Parameters

IF EXISTS

Do not throw an error if the schema does not exist. A notice is issued in this case.

name

The name of a schema.

CASCADE

Automatically drop objects (tables, functions, etc.) that are contained in the schema, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the schema if it contains any objects. This is the default.

Notes

Using the CASCADE option might make the command remove objects in other schemas besides the one(s) named.

Examples

To remove schema `mystuff` from the database, along with everything it contains:

```
DROP SCHEMA mystuff CASCADE;
```

Compatibility

DROP SCHEMA is fully conforming with the SQL standard, except that the standard only allows one schema to be dropped per command, and apart from the IF EXISTS option, which is a Postgres Pro extension.

See Also

[ALTER SCHEMA](#), [CREATE SCHEMA](#)

DROP SEQUENCE

DROP SEQUENCE — remove a sequence

Synopsis

```
DROP SEQUENCE [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

`DROP SEQUENCE` removes sequence number generators. A sequence can only be dropped by its owner or a superuser.

Parameters

`IF EXISTS`

Do not throw an error if the sequence does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of a sequence.

`CASCADE`

Automatically drop objects that depend on the sequence, and in turn all objects that depend on those objects (see [Section 5.15](#)).

`RESTRICT`

Refuse to drop the sequence if any objects depend on it. This is the default.

Examples

To remove the sequence `serial`:

```
DROP SEQUENCE serial;
```

Compatibility

`DROP SEQUENCE` conforms to the SQL standard, except that the standard only allows one sequence to be dropped per command, and apart from the `IF EXISTS` option, which is a Postgres Pro extension.

See Also

[CREATE SEQUENCE](#), [ALTER SEQUENCE](#)

DROP SERVER

DROP SERVER — remove a foreign server descriptor

Synopsis

```
DROP SERVER [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP SERVER removes an existing foreign server descriptor. To execute this command, the current user must be the owner of the server.

Parameters

IF EXISTS

Do not throw an error if the server does not exist. A notice is issued in this case.

name

The name of an existing server.

CASCADE

Automatically drop objects that depend on the server (such as user mappings), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the server if any objects depend on it. This is the default.

Examples

Drop a server `foo` if it exists:

```
DROP SERVER IF EXISTS foo;
```

Compatibility

DROP SERVER conforms to ISO/IEC 9075-9 (SQL/MED). The IF EXISTS clause is a Postgres Pro extension.

See Also

[CREATE SERVER](#), [ALTER SERVER](#)

DROP STATISTICS

DROP STATISTICS — remove extended statistics

Synopsis

```
DROP STATISTICS [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

`DROP STATISTICS` removes statistics object(s) from the database. Only the statistics object's owner, the schema owner, or a superuser can drop a statistics object.

Parameters

`IF EXISTS`

Do not throw an error if the statistics object does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of the statistics object to drop.

`CASCADE`

`RESTRICT`

These key words do not have any effect, since there are no dependencies on statistics.

Examples

To destroy two statistics objects in different schemas, without failing if they don't exist:

```
DROP STATISTICS IF EXISTS
    accounting.users_uid_creation,
    public.grants_user_role;
```

Compatibility

There is no `DROP STATISTICS` command in the SQL standard.

See Also

[ALTER STATISTICS](#), [CREATE STATISTICS](#)

DROP SUBSCRIPTION

DROP SUBSCRIPTION — remove a subscription

Synopsis

```
DROP SUBSCRIPTION [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP SUBSCRIPTION removes a subscription from the database cluster.

To execute this command the user must be the owner of the subscription.

DROP SUBSCRIPTION cannot be executed inside a transaction block if the subscription is associated with a replication slot. (You can use [ALTER SUBSCRIPTION](#) to unset the slot.)

Parameters

name

The name of a subscription to be dropped.

CASCADE

RESTRICT

These key words do not have any effect, since there are no dependencies on subscriptions.

Notes

When dropping a subscription that is associated with a replication slot on the remote host (the normal state), DROP SUBSCRIPTION will connect to the remote host and try to drop the replication slot (and any remaining table synchronization slots) as part of its operation. This is necessary so that the resources allocated for the subscription on the remote host are released. If this fails, either because the remote host is not reachable or because the remote replication slot cannot be dropped or does not exist or never existed, the DROP SUBSCRIPTION command will fail. To proceed in this situation, first disable the subscription by executing [ALTER SUBSCRIPTION ... DISABLE](#), and then disassociate it from the replication slot by executing [ALTER SUBSCRIPTION ... SET \(slot_name = NONE\)](#). After that, DROP SUBSCRIPTION will no longer attempt any actions on a remote host. Note that if the remote replication slot still exists, it (and any related table synchronization slots) should then be dropped manually; otherwise it/they will continue to reserve WAL and might eventually cause the disk to fill up. See also [Section 29.2.1](#).

If a subscription is associated with a replication slot, then DROP SUBSCRIPTION cannot be executed inside a transaction block.

Examples

Drop a subscription:

```
DROP SUBSCRIPTION mysub;
```

Compatibility

DROP SUBSCRIPTION is a Postgres Pro extension.

See Also

[CREATE SUBSCRIPTION](#), [ALTER SUBSCRIPTION](#)

DROP TABLE

DROP TABLE — remove a table

Synopsis

```
DROP TABLE [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP TABLE removes tables from the database. Only the table owner, the schema owner, and superuser can drop a table. To empty a table of rows without destroying the table, use [DELETE](#) or [TRUNCATE](#).

DROP TABLE always removes any indexes, rules, triggers, and constraints that exist for the target table. However, to drop a table that is referenced by a view or a foreign-key constraint of another table, CASCADE must be specified. (CASCADE will remove a dependent view entirely, but in the foreign-key case it will only remove the foreign-key constraint, not the other table entirely.)

Parameters

IF EXISTS

Do not throw an error if the table does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of the table to drop.

CASCADE

Automatically drop objects that depend on the table (such as views), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the table if any objects depend on it. This is the default.

Examples

To destroy two tables, `films` and `distributors`:

```
DROP TABLE films, distributors;
```

Compatibility

This command conforms to the SQL standard, except that the standard only allows one table to be dropped per command, and apart from the IF EXISTS option, which is a Postgres Pro extension.

See Also

[ALTER TABLE](#), [CREATE TABLE](#)

DROP TABLESPACE

DROP TABLESPACE — remove a tablespace

Synopsis

```
DROP TABLESPACE [ IF EXISTS ] name
```

Description

DROP TABLESPACE removes a tablespace from the system.

A tablespace can only be dropped by its owner or a superuser. The tablespace must be empty of all database objects before it can be dropped. It is possible that objects in other databases might still reside in the tablespace even if no objects in the current database are using the tablespace. Also, if the tablespace is listed in the [temp_tablespaces](#) setting of any active session, the DROP might fail due to temporary files residing in the tablespace.

Parameters

IF EXISTS

Do not throw an error if the tablespace does not exist. A notice is issued in this case.

name

The name of a tablespace.

Notes

DROP TABLESPACE cannot be executed inside a transaction block.

Examples

To remove tablespace `mystuff` from the system:

```
DROP TABLESPACE mystuff;
```

Compatibility

DROP TABLESPACE is a Postgres Pro extension.

See Also

[CREATE TABLESPACE](#), [ALTER TABLESPACE](#)

DROP TEXT SEARCH CONFIGURATION

DROP TEXT SEARCH CONFIGURATION — remove a text search configuration

Synopsis

```
DROP TEXT SEARCH CONFIGURATION [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP TEXT SEARCH CONFIGURATION **drops** an existing text search configuration. To execute this command you must be the owner of the configuration.

Parameters

IF EXISTS

Do not throw an error if the text search configuration does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing text search configuration.

CASCADE

Automatically drop objects that depend on the text search configuration, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the text search configuration if any objects depend on it. This is the default.

Examples

Remove the text search configuration `my_english`:

```
DROP TEXT SEARCH CONFIGURATION my_english;
```

This command will not succeed if there are any existing indexes that reference the configuration in `to_tsvector` calls. Add `CASCADE` to drop such indexes along with the text search configuration.

Compatibility

There is no DROP TEXT SEARCH CONFIGURATION statement in the SQL standard.

See Also

[ALTER TEXT SEARCH CONFIGURATION](#), [CREATE TEXT SEARCH CONFIGURATION](#)

DROP TEXT SEARCH DICTIONARY

DROP TEXT SEARCH DICTIONARY — remove a text search dictionary

Synopsis

```
DROP TEXT SEARCH DICTIONARY [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP TEXT SEARCH DICTIONARY drops an existing text search dictionary. To execute this command you must be the owner of the dictionary.

Parameters

IF EXISTS

Do not throw an error if the text search dictionary does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing text search dictionary.

CASCADE

Automatically drop objects that depend on the text search dictionary, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the text search dictionary if any objects depend on it. This is the default.

Examples

Remove the text search dictionary `english`:

```
DROP TEXT SEARCH DICTIONARY english;
```

This command will not succeed if there are any existing text search configurations that use the dictionary. Add `CASCADE` to drop such configurations along with the dictionary.

Compatibility

There is no DROP TEXT SEARCH DICTIONARY statement in the SQL standard.

See Also

[ALTER TEXT SEARCH DICTIONARY](#), [CREATE TEXT SEARCH DICTIONARY](#)

DROP TEXT SEARCH PARSER

DROP TEXT SEARCH PARSER — remove a text search parser

Synopsis

```
DROP TEXT SEARCH PARSER [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP TEXT SEARCH PARSER drops an existing text search parser. You must be a superuser to use this command.

Parameters

IF EXISTS

Do not throw an error if the text search parser does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing text search parser.

CASCADE

Automatically drop objects that depend on the text search parser, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the text search parser if any objects depend on it. This is the default.

Examples

Remove the text search parser `my_parser`:

```
DROP TEXT SEARCH PARSER my_parser;
```

This command will not succeed if there are any existing text search configurations that use the parser. Add `CASCADE` to drop such configurations along with the parser.

Compatibility

There is no DROP TEXT SEARCH PARSER statement in the SQL standard.

See Also

[ALTER TEXT SEARCH PARSER](#), [CREATE TEXT SEARCH PARSER](#)

DROP TEXT SEARCH TEMPLATE

DROP TEXT SEARCH TEMPLATE — remove a text search template

Synopsis

```
DROP TEXT SEARCH TEMPLATE [ IF EXISTS ] name [ CASCADE | RESTRICT ]
```

Description

DROP TEXT SEARCH TEMPLATE drops an existing text search template. You must be a superuser to use this command.

Parameters

IF EXISTS

Do not throw an error if the text search template does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of an existing text search template.

CASCADE

Automatically drop objects that depend on the text search template, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the text search template if any objects depend on it. This is the default.

Examples

Remove the text search template `thesaurus`:

```
DROP TEXT SEARCH TEMPLATE thesaurus;
```

This command will not succeed if there are any existing text search dictionaries that use the template. Add `CASCADE` to drop such dictionaries along with the template.

Compatibility

There is no DROP TEXT SEARCH TEMPLATE statement in the SQL standard.

See Also

[ALTER TEXT SEARCH TEMPLATE](#), [CREATE TEXT SEARCH TEMPLATE](#)

DROP TRANSFORM

DROP TRANSFORM — remove a transform

Synopsis

```
DROP TRANSFORM [ IF EXISTS ] FOR type_name LANGUAGE lang_name [ CASCADE | RESTRICT ]
```

Description

DROP TRANSFORM removes a previously defined transform.

To be able to drop a transform, you must own the type and the language. These are the same privileges that are required to create a transform.

Parameters

IF EXISTS

Do not throw an error if the transform does not exist. A notice is issued in this case.

type_name

The name of the data type of the transform.

lang_name

The name of the language of the transform.

CASCADE

Automatically drop objects that depend on the transform, and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the transform if any objects depend on it. This is the default.

Examples

To drop the transform for type `hstore` and language `plpython3u`:

```
DROP TRANSFORM FOR hstore LANGUAGE plpython3u;
```

Compatibility

This form of DROP TRANSFORM is a Postgres Pro extension. See [CREATE TRANSFORM](#) for details.

See Also

[CREATE TRANSFORM](#)

DROP TRIGGER

DROP TRIGGER — remove a trigger

Synopsis

```
DROP TRIGGER [ IF EXISTS ] name ON table_name [ CASCADE | RESTRICT ]
```

Description

`DROP TRIGGER` removes an existing trigger definition. To execute this command, the current user must be the owner of the table for which the trigger is defined.

Parameters

`IF EXISTS`

Do not throw an error if the trigger does not exist. A notice is issued in this case.

name

The name of the trigger to remove.

table_name

The name (optionally schema-qualified) of the table for which the trigger is defined.

`CASCADE`

Automatically drop objects that depend on the trigger, and in turn all objects that depend on those objects (see [Section 5.15](#)).

`RESTRICT`

Refuse to drop the trigger if any objects depend on it. This is the default.

Examples

Destroy the trigger `if_dist_exists` on the table `films`:

```
DROP TRIGGER if_dist_exists ON films;
```

Compatibility

The `DROP TRIGGER` statement in Postgres Pro is incompatible with the SQL standard. In the SQL standard, trigger names are not local to tables, so the command is simply `DROP TRIGGER name`.

See Also

[CREATE TRIGGER](#)

DROP TYPE

DROP TYPE — remove a data type

Synopsis

```
DROP TYPE [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

DROP TYPE removes a user-defined data type. Only the owner of a type can remove it.

Parameters

IF EXISTS

Do not throw an error if the type does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of the data type to remove.

CASCADE

Automatically drop objects that depend on the type (such as table columns, functions, and operators), and in turn all objects that depend on those objects (see [Section 5.15](#)).

RESTRICT

Refuse to drop the type if any objects depend on it. This is the default.

Examples

To remove the data type `box`:

```
DROP TYPE box;
```

Compatibility

This command is similar to the corresponding command in the SQL standard, apart from the `IF EXISTS` option, which is a Postgres Pro extension. But note that much of the `CREATE TYPE` command and the data type extension mechanisms in Postgres Pro differ from the SQL standard.

See Also

[ALTER TYPE](#), [CREATE TYPE](#)

DROP USER

DROP USER — remove a database role

Synopsis

```
DROP USER [ IF EXISTS ] name [, ...]
```

Description

DROP USER is simply an alternate spelling of [DROP ROLE](#).

Compatibility

The DROP USER statement is a Postgres Pro extension. The SQL standard leaves the definition of users to the implementation.

See Also

[DROP ROLE](#)

DROP VIEW

DROP VIEW — remove a view

Synopsis

```
DROP VIEW [ IF EXISTS ] name [, ...] [ CASCADE | RESTRICT ]
```

Description

`DROP VIEW` drops an existing view. To execute this command you must be the owner of the view.

Parameters

`IF EXISTS`

Do not throw an error if the view does not exist. A notice is issued in this case.

name

The name (optionally schema-qualified) of the view to remove.

`CASCADE`

Automatically drop objects that depend on the view (such as other views), and in turn all objects that depend on those objects (see [Section 5.15](#)).

`RESTRICT`

Refuse to drop the view if any objects depend on it. This is the default.

Examples

This command will remove the view called `kinds`:

```
DROP VIEW kinds;
```

Compatibility

This command conforms to the SQL standard, except that the standard only allows one view to be dropped per command, and apart from the `IF EXISTS` option, which is a Postgres Pro extension.

See Also

[ALTER VIEW](#), [CREATE VIEW](#)

END

END — commit the current transaction

Synopsis

```
END [ AUTONOMOUS ] [ WORK | TRANSACTION ] [ AND [ NO ] CHAIN ]
```

Description

END commits the current transaction. All changes made by the transaction become visible to others and are guaranteed to be durable if a crash occurs. This command is a Postgres Pro extension that is equivalent to [COMMIT](#).

Parameters

AUTONOMOUS

Optional key word that can be used when committing an autonomous transaction. For details on autonomous transactions, see [Chapter 16](#).

WORK

TRANSACTION

AUTONOMOUS

Optional key words. They have no effect.

AND CHAIN

If AND CHAIN is specified, a new transaction is immediately started with the same transaction characteristics (see [SET TRANSACTION](#)) as the just finished one. Otherwise, no new transaction is started.

Notes

Use [ROLLBACK](#) to abort a transaction.

Issuing END when not inside a transaction does no harm, but it will provoke a warning message.

Examples

To commit the current transaction and make all changes permanent:

```
END;
```

Compatibility

END is a Postgres Pro extension that provides functionality equivalent to [COMMIT](#), which is specified in the SQL standard.

See Also

[BEGIN](#), [COMMIT](#), [ROLLBACK](#)

EXECUTE

EXECUTE — execute a prepared statement

Synopsis

```
EXECUTE name [ ( parameter [, ...] ) ]
```

Description

EXECUTE is used to execute a previously prepared statement. Since prepared statements only exist for the duration of a session, the prepared statement must have been created by a `PREPARE` statement executed earlier in the current session.

If the `PREPARE` statement that created the statement specified some parameters, a compatible set of parameters must be passed to the `EXECUTE` statement, or else an error is raised. Note that (unlike functions) prepared statements are not overloaded based on the type or number of their parameters; the name of a prepared statement must be unique within a database session.

For more information on the creation and usage of prepared statements, see [PREPARE](#).

Parameters

name

The name of the prepared statement to execute.

parameter

The actual value of a parameter to the prepared statement. This must be an expression yielding a value that is compatible with the data type of this parameter, as was determined when the prepared statement was created.

Outputs

The command tag returned by `EXECUTE` is that of the prepared statement, and not `EXECUTE`.

Examples

Examples are given in [Examples](#) in the [PREPARE](#) documentation.

Compatibility

The SQL standard includes an `EXECUTE` statement, but it is only for use in embedded SQL. This version of the `EXECUTE` statement also uses a somewhat different syntax.

See Also

[DEALLOCATE](#), [PREPARE](#)

EXPLAIN

EXPLAIN — show the execution plan of a statement

Synopsis

```
EXPLAIN [ ( option [, ...] ) ] statement
```

where *option* can be one of:

```
ANALYZE [ boolean ]
VERBOSE [ boolean ]
COSTS [ boolean ]
SETTINGS [ boolean ]
GENERIC_PLAN [ boolean ]
BUFFERS [ boolean ]
SERIALIZE [ { NONE | TEXT | BINARY } ]
WAL [ boolean ]
TIMING [ boolean ]
SUMMARY [ boolean ]
MEMORY [ boolean ]
FORMAT { TEXT | XML | JSON | YAML }
```

Description

This command displays the execution plan that the Postgres Pro planner generates for the supplied statement. The execution plan shows how the table(s) referenced by the statement will be scanned — by plain sequential scan, index scan, etc. — and if multiple tables are referenced, what join algorithms will be used to bring together the required rows from each input table.

The most critical part of the display is the estimated statement execution cost, which is the planner's guess at how long it will take to run the statement (measured in cost units that are arbitrary, but conventionally mean disk page fetches). Actually two numbers are shown: the start-up cost before the first row can be returned, and the total cost to return all the rows. For most queries the total cost is what matters, but in contexts such as a subquery in `EXISTS`, the planner will choose the smallest start-up cost instead of the smallest total cost (since the executor will stop after getting one row, anyway). Also, if you limit the number of rows to return with a `LIMIT` clause, the planner makes an appropriate interpolation between the endpoint costs to estimate which plan is really the cheapest.

The `ANALYZE` option causes the statement to be actually executed, not only planned. Then actual run time statistics are added to the display, including the total elapsed time expended within each plan node (in milliseconds) and the total number of rows it actually returned. This is useful for seeing whether the planner's estimates are close to reality.

Important

Keep in mind that the statement is actually executed when the `ANALYZE` option is used. Although `EXPLAIN` will discard any output that a `SELECT` would return, other side effects of the statement will happen as usual. If you wish to use `EXPLAIN ANALYZE` on an `INSERT`, `UPDATE`, `DELETE`, `MERGE`, `CREATE TABLE AS`, or `EXECUTE` statement without letting the command affect your data, use this approach:

```
BEGIN;
EXPLAIN ANALYZE ...;
ROLLBACK;
```

Parameters

ANALYZE

Carry out the command and show actual run times and other statistics. This parameter defaults to `FALSE`.

VERBOSE

Display additional information regarding the plan. Specifically, include the output column list for each node in the plan tree, schema-qualify table and function names, always label variables in expressions with their range table alias, and always print the name of each trigger for which statistics are displayed. The query identifier will also be displayed if one has been computed, see [compute_query_id](#) for more details. This parameter defaults to `FALSE`.

COSTS

Include information on the estimated startup and total cost of each plan node, as well as the estimated number of rows and the estimated width of each row. This parameter defaults to `TRUE`.

SETTINGS

Include information on configuration parameters. Specifically, include options affecting query planning with value different from the built-in default value. This parameter defaults to `FALSE`.

GENERIC_PLAN

Allow the statement to contain parameter placeholders like `$1`, and generate a generic plan that does not depend on the values of those parameters. See [PREPARE](#) for details about generic plans and the types of statement that support parameters. This parameter cannot be used together with `ANALYZE`. It defaults to `FALSE`.

BUFFERS

Include information on buffer usage. Specifically, include the number of shared blocks hit, read, dirtied, and written, the number of local blocks hit, read, dirtied, and written, the number of temp blocks read and written, and the time spent reading and writing data file blocks, local blocks and temporary file blocks (in milliseconds) if [track_io_timing](#) is enabled. A *hit* means that a read was avoided because the block was found already in cache when needed. Shared blocks contain data from regular tables and indexes; local blocks contain data from temporary tables and indexes; while temporary blocks contain short-term working data used in sorts, hashes, Materialize plan nodes, and similar cases. The number of blocks *dirtied* indicates the number of previously unmodified blocks that were changed by this query; while the number of blocks *written* indicates the number of previously-dirtied blocks evicted from cache by this backend during query processing. The number of blocks shown for an upper-level node includes those used by all its child nodes. In text format, only non-zero values are printed. This parameter defaults to `FALSE`.

SERIALIZE

Include information on the cost of *serializing* the query's output data, that is converting it to text or binary format to send to the client. This can be a significant part of the time required for regular execution of the query, if the datatype output functions are expensive or if TOASTed values must be fetched from out-of-line storage. `EXPLAIN`'s default behavior, `SERIALIZE NONE`, does not perform these conversions. If `SERIALIZE TEXT` or `SERIALIZE BINARY` is specified, the appropriate conversions are performed, and the time spent doing so is measured (unless `TIMING OFF` is specified). If the `BUFFERS` option is also specified, then any buffer accesses involved in the conversions are counted too. In no case, however, will `EXPLAIN` actually send the resulting data to the client; hence network transmission costs cannot be investigated this way. Serialization may only be enabled when `ANALYZE` is also enabled. If `SERIALIZE` is written without an argument, `TEXT` is assumed.

WAL

Include information on WAL record generation. Specifically, include the number of records, number of full page images (fpi) and the amount of WAL generated in bytes. In text format, only non-zero

values are printed. This parameter may only be used when `ANALYZE` is also enabled. It defaults to `FALSE`.

TIMING

Include actual startup time and time spent in each node in the output. The overhead of repeatedly reading the system clock can slow down the query significantly on some systems, so it may be useful to set this parameter to `FALSE` when only actual row counts, and not exact times, are needed. Run time of the entire statement is always measured, even when node-level timing is turned off with this option. This parameter may only be used when `ANALYZE` is also enabled. It defaults to `TRUE`.

SUMMARY

Include summary information (e.g., totaled timing information) after the query plan. Summary information is included by default when `ANALYZE` is used but otherwise is not included by default, but can be enabled using this option. Planning time in `EXPLAIN EXECUTE` includes the time required to fetch the plan from the cache and the time required for re-planning, if necessary.

MEMORY

Include information on memory consumption by the query planning phase. Specifically, include the precise amount of storage used by planner in-memory structures, as well as total memory considering allocation overhead. This parameter defaults to `FALSE`.

FORMAT

Specify the output format, which can be `TEXT`, `XML`, `JSON`, or `YAML`. Non-text output contains the same information as the text output format, but is easier for programs to parse. This parameter defaults to `TEXT`.

boolean

Specifies whether the selected option should be turned on or off. You can write `TRUE`, `ON`, or `1` to enable the option, and `FALSE`, `OFF`, or `0` to disable it. The *boolean* value can also be omitted, in which case `TRUE` is assumed.

statement

Any `SELECT`, `INSERT`, `UPDATE`, `DELETE`, `MERGE`, `VALUES`, `EXECUTE`, `DECLARE`, `CREATE TABLE AS`, or `CREATE MATERIALIZED VIEW AS statement`, whose execution plan you wish to see.

Outputs

The command's result is a textual description of the plan selected for the *statement*, optionally annotated with execution statistics. [Section 14.1](#) describes the information provided.

Notes

In order to allow the Postgres Pro query planner to make reasonably informed decisions when optimizing queries, the `pg_statistic` data should be up-to-date for all tables used in the query. Normally the [auto-vacuum daemon](#) will take care of that automatically. But if a table has recently had substantial changes in its contents, you might need to do a manual `ANALYZE` rather than wait for autovacuum to catch up with the changes.

In order to measure the run-time cost of each node in the execution plan, the current implementation of `EXPLAIN ANALYZE` adds profiling overhead to query execution. As a result, running `EXPLAIN ANALYZE` on a query can sometimes take significantly longer than executing the query normally. The amount of overhead depends on the nature of the query, as well as the platform being used. The worst case occurs for plan nodes that in themselves require very little time per execution, and on machines that have relatively slow operating system calls for obtaining the time of day.

Examples

To show the plan for a simple query on a table with a single `integer` column and 10000 rows:

```
EXPLAIN SELECT * FROM foo;
```

QUERY PLAN

```
-----
Seq Scan on foo  (cost=0.00..155.00 rows=10000 width=4)
(1 row)
```

Here is the same query, with JSON output formatting:

```
EXPLAIN (FORMAT JSON) SELECT * FROM foo;
QUERY PLAN
```

```
-----
[
  {
    "Plan": {
      "Node Type": "Seq Scan",
      "Relation Name": "foo",
      "Alias": "foo",
      "Startup Cost": 0.00,
      "Total Cost": 155.00,
      "Plan Rows": 10000,
      "Plan Width": 4
    }
  }
]
(1 row)
```

If there is an index and we use a query with an indexable WHERE condition, EXPLAIN might show a different plan:

```
EXPLAIN SELECT * FROM foo WHERE i = 4;
```

QUERY PLAN

```
-----
Index Scan using fi on foo  (cost=0.00..5.98 rows=1 width=4)
Index Cond: (i = 4)
(2 rows)
```

Here is the same query, but in YAML format:

```
EXPLAIN (FORMAT YAML) SELECT * FROM foo WHERE i='4';
QUERY PLAN
```

```
-----
- Plan:
  Node Type: "Index Scan"
  Scan Direction: "Forward"
  Index Name: "fi"
  Relation Name: "foo"
  Alias: "foo"
  Startup Cost: 0.00
  Total Cost: 5.98
  Plan Rows: 1
  Plan Width: 4
  Index Cond: "(i = 4)"
(1 row)
```

XML format is left as an exercise for the reader.

Here is the same plan with cost estimates suppressed:

```
EXPLAIN (COSTS FALSE) SELECT * FROM foo WHERE i = 4;
```

QUERY PLAN

```
-----
Index Scan using fi on foo
  Index Cond: (i = 4)
(2 rows)
```

Here is an example of a query plan for a query using an aggregate function:

```
EXPLAIN SELECT sum(i) FROM foo WHERE i < 10;
```

QUERY PLAN

```
-----
Aggregate  (cost=23.93..23.93 rows=1 width=4)
->  Index Scan using fi on foo  (cost=0.00..23.92 rows=6 width=4)
     Index Cond: (i < 10)
(3 rows)
```

Here is an example of using `EXPLAIN EXECUTE` to display the execution plan for a prepared query:

```
PREPARE query(int, int) AS SELECT sum(bar) FROM test
  WHERE id > $1 AND id < $2
  GROUP BY foo;
```

```
EXPLAIN ANALYZE EXECUTE query(100, 200);
```

QUERY PLAN

```
-----
HashAggregate  (cost=10.77..10.87 rows=10 width=12) (actual time=0.043..0.044 rows=10
loops=1)
  Group Key: foo
  Batches: 1  Memory Usage: 24kB
->  Index Scan using test_pkey on test  (cost=0.29..10.27 rows=99 width=8) (actual
time=0.009..0.025 rows=99 loops=1)
     Index Cond: ((id > 100) AND (id < 200))
Planning Time: 0.244 ms
Execution Time: 0.073 ms
(7 rows)
```

Of course, the specific numbers shown here depend on the actual contents of the tables involved. Also note that the numbers, and even the selected query strategy, might vary between Postgres Pro releases due to planner improvements. In addition, the `ANALYZE` command uses random sampling to estimate data statistics; therefore, it is possible for cost estimates to change after a fresh run of `ANALYZE`, even if the actual distribution of data in the table has not changed.

Notice that the previous example showed a “custom” plan for the specific parameter values given in `EXECUTE`. We might also wish to see the generic plan for a parameterized query, which can be done with `GENERIC_PLAN`:

```
EXPLAIN (GENERIC_PLAN)
  SELECT sum(bar) FROM test
  WHERE id > $1 AND id < $2
  GROUP BY foo;
```

QUERY PLAN

```
-----
HashAggregate  (cost=26.79..26.89 rows=10 width=12)
  Group Key: foo
->  Index Scan using test_pkey on test  (cost=0.29..24.29 rows=500 width=8)
```

```
      Index Cond: ((id > $1) AND (id < $2))
(4 rows)
```

In this case the parser correctly inferred that `$1` and `$2` should have the same data type as `id`, so the lack of parameter type information from `PREPARE` was not a problem. In other cases it might be necessary to explicitly specify types for the parameter symbols, which can be done by casting them, for example:

```
EXPLAIN (GENERIC_PLAN)
  SELECT sum(bar) FROM test
     WHERE id > $1::integer AND id < $2::integer
     GROUP BY foo;
```

Compatibility

There is no `EXPLAIN` statement defined in the SQL standard.

The following syntax was used before Postgres Pro version 9.0 and is still supported:

```
EXPLAIN [ ANALYZE ] [ VERBOSE ] statement
```

Note that in this syntax, the options must be specified in exactly the order shown.

See Also

[ANALYZE](#)

FETCH

FETCH — retrieve rows from a query using a cursor

Synopsis

```
FETCH [ direction ] [ FROM | IN ] cursor_name
```

where *direction* can be one of:

```
NEXT  
PRIOR  
FIRST  
LAST  
ABSOLUTE count  
RELATIVE count  
count  
ALL  
FORWARD  
FORWARD count  
FORWARD ALL  
BACKWARD  
BACKWARD count  
BACKWARD ALL
```

Description

FETCH retrieves rows using a previously-created cursor.

A cursor has an associated position, which is used by FETCH. The cursor position can be before the first row of the query result, on any particular row of the result, or after the last row of the result. When created, a cursor is positioned before the first row. After fetching some rows, the cursor is positioned on the row most recently retrieved. If FETCH runs off the end of the available rows then the cursor is left positioned after the last row, or before the first row if fetching backward. FETCH ALL or FETCH BACKWARD ALL will always leave the cursor positioned after the last row or before the first row.

The forms NEXT, PRIOR, FIRST, LAST, ABSOLUTE, RELATIVE fetch a single row after moving the cursor appropriately. If there is no such row, an empty result is returned, and the cursor is left positioned before the first row or after the last row as appropriate.

The forms using FORWARD and BACKWARD retrieve the indicated number of rows moving in the forward or backward direction, leaving the cursor positioned on the last-returned row (or after/before all rows, if the *count* exceeds the number of rows available).

RELATIVE 0, FORWARD 0, and BACKWARD 0 all request fetching the current row without moving the cursor, that is, re-fetching the most recently fetched row. This will succeed unless the cursor is positioned before the first row or after the last row; in which case, no row is returned.

Note

This page describes usage of cursors at the SQL command level. If you are trying to use cursors inside a PL/pgSQL function, the rules are different — see [Section 44.7.3](#).

Parameters

direction

direction defines the fetch direction and number of rows to fetch. It can be one of the following:

NEXT

Fetch the next row. This is the default if *direction* is omitted.

PRIOR

Fetch the prior row.

FIRST

Fetch the first row of the query (same as ABSOLUTE 1).

LAST

Fetch the last row of the query (same as ABSOLUTE -1).

ABSOLUTE *count*

Fetch the *count*'th row of the query, or the *abs(count)*'th row from the end if *count* is negative. Position before first row or after last row if *count* is out of range; in particular, ABSOLUTE 0 positions before the first row.

RELATIVE *count*

Fetch the *count*'th succeeding row, or the *abs(count)*'th prior row if *count* is negative. RELATIVE 0 re-fetches the current row, if any.

count

Fetch the next *count* rows (same as FORWARD *count*).

ALL

Fetch all remaining rows (same as FORWARD ALL).

FORWARD

Fetch the next row (same as NEXT).

FORWARD *count*

Fetch the next *count* rows. FORWARD 0 re-fetches the current row.

FORWARD ALL

Fetch all remaining rows.

BACKWARD

Fetch the prior row (same as PRIOR).

BACKWARD *count*

Fetch the prior *count* rows (scanning backwards). BACKWARD 0 re-fetches the current row.

BACKWARD ALL

Fetch all prior rows (scanning backwards).

count

count is a possibly-signed integer constant, determining the location or number of rows to fetch. For FORWARD and BACKWARD cases, specifying a negative *count* is equivalent to changing the sense of FORWARD and BACKWARD.

cursor_name

An open cursor's name.

Outputs

On successful completion, a `FETCH` command returns a command tag of the form

```
FETCH count
```

The *count* is the number of rows fetched (possibly zero). Note that in `psql`, the command tag will not actually be displayed, since `psql` displays the fetched rows instead.

Notes

The cursor should be declared with the `SCROLL` option if one intends to use any variants of `FETCH` other than `FETCH NEXT` or `FETCH FORWARD` with a positive count. For simple queries Postgres Pro will allow backwards fetch from cursors not declared with `SCROLL`, but this behavior is best not relied on. If the cursor is declared with `NO SCROLL`, no backward fetches are allowed.

`ABSOLUTE` fetches are not any faster than navigating to the desired row with a relative move: the underlying implementation must traverse all the intermediate rows anyway. Negative absolute fetches are even worse: the query must be read to the end to find the last row, and then traversed backward from there. However, rewinding to the start of the query (as with `FETCH ABSOLUTE 0`) is fast.

`DECLARE` is used to define a cursor. Use `MOVE` to change cursor position without retrieving data.

Examples

The following example traverses a table using a cursor:

```
BEGIN WORK;

-- Set up a cursor:
DECLARE liahona SCROLL CURSOR FOR SELECT * FROM films;

-- Fetch the first 5 rows in the cursor liahona:
FETCH FORWARD 5 FROM liahona;
```

code	title	did	date_prod	kind	len
BL101	The Third Man	101	1949-12-23	Drama	01:44
BL102	The African Queen	101	1951-08-11	Romantic	01:43
JL201	Une Femme est une Femme	102	1961-03-12	Romantic	01:25
P_301	Vertigo	103	1958-11-14	Action	02:08
P_302	Becket	103	1964-02-03	Drama	02:28

```
-- Fetch the previous row:
FETCH PRIOR FROM liahona;
```

code	title	did	date_prod	kind	len
P_301	Vertigo	103	1958-11-14	Action	02:08

```
-- Close the cursor and end the transaction:
CLOSE liahona;
COMMIT WORK;
```

Compatibility

The SQL standard defines `FETCH` for use in embedded SQL only. The variant of `FETCH` described here returns the data as if it were a `SELECT` result rather than placing it in host variables. Other than this point, `FETCH` is fully upward-compatible with the SQL standard.

The `FETCH` forms involving `FORWARD` and `BACKWARD`, as well as the forms `FETCH count` and `FETCH ALL`, in which `FORWARD` is implicit, are Postgres Pro extensions.

The SQL standard allows only `FROM` preceding the cursor name; the option to use `IN`, or to leave them out altogether, is an extension.

See Also

[CLOSE](#), [DECLARE](#), [MOVE](#)

GRANT

GRANT — define access privileges

Synopsis

```
GRANT { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER |
        MAINTAIN }
        [, ...] | ALL [ PRIVILEGES ] }
ON { [ TABLE ] table_name [, ...]
     | ALL TABLES IN SCHEMA schema_name [, ...] }
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { { SELECT | INSERT | UPDATE | REFERENCES } ( column_name [, ...] )
        [, ...] | ALL [ PRIVILEGES ] ( column_name [, ...] ) }
ON [ TABLE ] table_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { { USAGE | SELECT | UPDATE }
        [, ...] | ALL [ PRIVILEGES ] }
ON { SEQUENCE sequence_name [, ...]
     | ALL SEQUENCES IN SCHEMA schema_name [, ...] }
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { { CREATE | CONNECT | TEMPORARY | TEMP } [, ...] | ALL [ PRIVILEGES ] }
ON DATABASE database_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { USAGE | ALL [ PRIVILEGES ] }
ON DOMAIN domain_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { USAGE | ALL [ PRIVILEGES ] }
ON FOREIGN DATA WRAPPER fdw_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { USAGE | ALL [ PRIVILEGES ] }
ON FOREIGN SERVER server_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { EXECUTE | ALL [ PRIVILEGES ] }
ON { { FUNCTION | PROCEDURE | ROUTINE } routine_name [ ( [ [ argmode ] [ arg_name
] arg_type [, ...] ) ] ) [, ...]
     | ALL { FUNCTIONS | PROCEDURES | ROUTINES } IN SCHEMA schema_name [, ...] }
TO role_specification [, ...] [ WITH GRANT OPTION ]
[ GRANTED BY role_specification ]

GRANT { USAGE | ALL [ PRIVILEGES ] }
ON LANGUAGE lang_name [, ...]
TO role_specification [, ...] [ WITH GRANT OPTION ]
```

```
[ GRANTED BY role_specification ]

GRANT { { SELECT | UPDATE } [, ...] | ALL [ PRIVILEGES ] }
  ON LARGE OBJECT loid [, ...]
  TO role_specification [, ...] [ WITH GRANT OPTION ]
  [ GRANTED BY role_specification ]

GRANT { { SET | ALTER SYSTEM } [, ...] | ALL [ PRIVILEGES ] }
  ON PARAMETER configuration_parameter [, ...]
  TO role_specification [, ...] [ WITH GRANT OPTION ]
  [ GRANTED BY role_specification ]

GRANT { { CREATE | USAGE } [, ...] | ALL [ PRIVILEGES ] }
  ON SCHEMA schema_name [, ...]
  TO role_specification [, ...] [ WITH GRANT OPTION ]
  [ GRANTED BY role_specification ]

GRANT { CREATE | ALL [ PRIVILEGES ] }
  ON TABLESPACE tablespace_name [, ...]
  TO role_specification [, ...] [ WITH GRANT OPTION ]
  [ GRANTED BY role_specification ]

GRANT { USAGE | ALL [ PRIVILEGES ] }
  ON TYPE type_name [, ...]
  TO role_specification [, ...] [ WITH GRANT OPTION ]
  [ GRANTED BY role_specification ]

GRANT role_name [, ...] TO role_specification [, ...]
  [ WITH { ADMIN | INHERIT | SET } { OPTION | TRUE | FALSE } ]
  [ GRANTED BY role_specification ]
```

where *role_specification* can be:

```
[ GROUP ] role_name
| PUBLIC
| CURRENT_ROLE
| CURRENT_USER
| SESSION_USER
```

Description

The GRANT command has two basic variants: one that grants privileges on a database object (table, column, view, foreign table, sequence, database, foreign-data wrapper, foreign server, function, procedure, procedural language, large object, configuration parameter, schema, tablespace, or type), and one that grants membership in a role. These variants are similar in many ways, but they are different enough to be described separately.

GRANT on Database Objects

This variant of the GRANT command gives specific privileges on a database object to one or more roles. These privileges are added to those already granted, if any.

The key word PUBLIC indicates that the privileges are to be granted to all roles, including those that might be created later. PUBLIC can be thought of as an implicitly defined group that always includes all roles. Any particular role will have the sum of privileges granted directly to it, privileges granted to any role it is presently a member of, and privileges granted to PUBLIC.

If WITH GRANT OPTION is specified, the recipient of the privilege can in turn grant it to others. Without a grant option, the recipient cannot do that. Grant options cannot be granted to PUBLIC.

If `GRANTED BY` is specified, the specified grantor must be the current user. This clause is currently present in this form only for SQL compatibility.

There is no need to grant privileges to the owner of an object (usually the user that created it), as the owner has all privileges by default. (The owner could, however, choose to revoke some of their own privileges for safety.)

The right to drop an object, or to alter its definition in any way, is not treated as a grantable privilege; it is inherent in the owner, and cannot be granted or revoked. (However, a similar effect can be obtained by granting or revoking membership in the role that owns the object; see below.) The owner implicitly has all grant options for the object, too.

The possible privileges are:

SELECT
INSERT
UPDATE
DELETE
TRUNCATE
REFERENCES
TRIGGER
CREATE
CONNECT
TEMPORARY
EXECUTE
USAGE
SET
ALTER SYSTEM
MAINTAIN

Specific types of privileges, as defined in [Section 5.8](#).

TEMP

Alternative spelling for `TEMPORARY`.

ALL PRIVILEGES

Grant all of the privileges available for the object's type. The `PRIVILEGES` key word is optional in Postgres Pro, though it is required by strict SQL.

The `FUNCTION` syntax works for plain functions, aggregate functions, and window functions, but not for procedures; use `PROCEDURE` for those. Alternatively, use `ROUTINE` to refer to a function, aggregate function, window function, or procedure regardless of its precise type.

There is also an option to grant privileges on all objects of the same type within one or more schemas. This functionality is currently supported only for tables, sequences, functions, and procedures. `ALL TABLES` also affects views and foreign tables, just like the specific-object `GRANT` command. `ALL FUNCTIONS` also affects aggregate and window functions, but not procedures, again just like the specific-object `GRANT` command. Use `ALL ROUTINES` to include procedures.

GRANT on Roles

This variant of the `GRANT` command grants membership in a role to one or more other roles, and the modification of membership options `SET`, `INHERIT`, and `ADMIN`; see [Section 21.3](#) for details. Membership in a role is significant because it potentially allows access to the privileges granted to a role to each of its members, and potentially also the ability to make changes to the role itself. However, the actual permissions conferred depend on the options associated with the grant. To modify that options of an existing membership, simply specify the membership with updated option values.

Each of the options described below can be set to either `TRUE` or `FALSE`. The keyword `OPTION` is accepted as a synonym for `TRUE`, so that `WITH ADMIN OPTION` is a synonym for `WITH ADMIN TRUE`. When altering an existing membership the omission of an option results in the current value being retained.

The `ADMIN` option allows the member to in turn grant membership in the role to others, and revoke membership in the role as well. Without the admin option, ordinary users cannot do that. A role is not considered to hold `WITH ADMIN OPTION` on itself. Database superusers can grant or revoke membership in any role to anyone. This option defaults to `FALSE`.

The `INHERIT` option controls the inheritance status of the new membership; see [Section 21.3](#) for details on inheritance. If it is set to `TRUE`, it causes the new member to inherit from the granted role. If set to `FALSE`, the new member does not inherit. If unspecified when creating a new role membership, this defaults to the inheritance attribute of the new member.

The `SET` option, if it is set to `TRUE`, allows the member to change to the granted role using the `SET ROLE` command. If a role is an indirect member of another role, it can use `SET ROLE` to change to that role only if there is a chain of grants each of which has `SET TRUE`. This option defaults to `TRUE`.

To create an object owned by another role or give ownership of an existing object to another role, you must have the ability to `SET ROLE` to that role; otherwise, commands such as `ALTER ... OWNER TO` or `CREATE DATABASE ... OWNER` will fail. However, a user who inherits the privileges of a role but does not have the ability to `SET ROLE` to that role may be able to obtain full access to the role by manipulating existing objects owned by that role (e.g. they could redefine an existing function to act as a Trojan horse). Therefore, if a role's privileges are to be inherited but should not be accessible via `SET ROLE`, it should not own any SQL objects.

If `GRANTED BY` is specified, the grant is recorded as having been done by the specified role. A user can only attribute a grant to another role if they possess the privileges of that role. The role recorded as the grantor must have `ADMIN OPTION` on the target role, unless it is the bootstrap superuser. When a grant is recorded as having a grantor other than the bootstrap superuser, it depends on the grantor continuing to possess `ADMIN OPTION` on the role; so, if `ADMIN OPTION` is revoked, dependent grants must be revoked as well.

Unlike the case with privileges, membership in a role cannot be granted to `PUBLIC`. Note also that this form of the command does not allow the noise word `GROUP` in *role_specification*.

Notes

The `REVOKE` command is used to revoke access privileges.

Since PostgreSQL 8.1, the concepts of users and groups have been unified into a single kind of entity called a role. It is therefore no longer necessary to use the keyword `GROUP` to identify whether a grantee is a user or a group. `GROUP` is still allowed in the command, but it is a noise word.

A user may perform `SELECT`, `INSERT`, etc. on a column if they hold that privilege for either the specific column or its whole table. Granting the privilege at the table level and then revoking it for one column will not do what one might wish: the table-level grant is unaffected by a column-level operation.

When a non-owner of an object attempts to `GRANT` privileges on the object, the command will fail outright if the user has no privileges whatsoever on the object. As long as some privilege is available, the command will proceed, but it will grant only those privileges for which the user has grant options. The `GRANT ALL PRIVILEGES` forms will issue a warning message if no grant options are held, while the other forms will issue a warning if grant options for any of the privileges specifically named in the command are not held. (In principle these statements apply to the object owner as well, but since the owner is always treated as holding all grant options, the cases can never occur.)

It should be noted that database superusers can access all objects regardless of object privilege settings. This is comparable to the rights of `root` in a Unix system. As with `root`, it's unwise to operate as a superuser except when absolutely necessary.

If a superuser chooses to issue a `GRANT` or `REVOKE` command, the command is performed as though it were issued by the owner of the affected object. In particular, privileges granted via such a command will appear to have been granted by the object owner. (For role membership, the membership appears to have been granted by the bootstrap superuser.)

`GRANT` and `REVOKE` can also be done by a role that is not the owner of the affected object, but is a member of the role that owns the object, or is a member of a role that holds privileges `WITH GRANT OPTION` on the object. In this case the privileges will be recorded as having been granted by the role that actually owns the object or holds the privileges `WITH GRANT OPTION`. For example, if table `t1` is owned by role `g1`, of which role `u1` is a member, then `u1` can grant privileges on `t1` to `u2`, but those privileges will appear to have been granted directly by `g1`. Any other member of role `g1` could revoke them later.

If the role executing `GRANT` holds the required privileges indirectly via more than one role membership path, it is unspecified which containing role will be recorded as having done the grant. In such cases it is best practice to use `SET ROLE` to become the specific role you want to do the `GRANT` as.

Granting permission on a table does not automatically extend permissions to any sequences used by the table, including sequences tied to `SERIAL` columns. Permissions on sequences must be set separately.

See [Section 5.8](#) for more information about specific privilege types, as well as how to inspect objects' privileges.

Examples

Grant insert privilege to all users on table `films`:

```
GRANT INSERT ON films TO PUBLIC;
```

Grant all available privileges to user `manuel` on view `kinds`:

```
GRANT ALL PRIVILEGES ON kinds TO manuel;
```

Note that while the above will indeed grant all privileges if executed by a superuser or the owner of `kinds`, when executed by someone else it will only grant those permissions for which the someone else has grant options.

Grant membership in role `admins` to user `joe`:

```
GRANT admins TO joe;
```

Compatibility

According to the SQL standard, the `PRIVILEGES` key word in `ALL PRIVILEGES` is required. The SQL standard does not support setting the privileges on more than one object per command.

Postgres Pro allows an object owner to revoke their own ordinary privileges: for example, a table owner can make the table read-only to themselves by revoking their own `INSERT`, `UPDATE`, `DELETE`, and `TRUNCATE` privileges. This is not possible according to the SQL standard. The reason is that Postgres Pro treats the owner's privileges as having been granted by the owner to themselves; therefore they can revoke them too. In the SQL standard, the owner's privileges are granted by an assumed entity “`_SYSTEM`”. Not being “`_SYSTEM`”, the owner cannot revoke these rights.

According to the SQL standard, grant options can be granted to `PUBLIC`; Postgres Pro only supports granting grant options to roles.

The SQL standard allows the `GRANTED BY` option to specify only `CURRENT_USER` or `CURRENT_ROLE`. The other variants are Postgres Pro extensions.

The SQL standard provides for a `USAGE` privilege on other kinds of objects: character sets, collations, translations.

In the SQL standard, sequences only have a `USAGE` privilege, which controls the use of the `NEXT VALUE FOR` expression, which is equivalent to the function `nextval` in Postgres Pro. The sequence privileges `SELECT` and `UPDATE` are Postgres Pro extensions. The application of the sequence `USAGE` privilege to the `currval` function is also a Postgres Pro extension (as is the function itself).

Privileges on databases, tablespaces, schemas, languages, and configuration parameters are Postgres Pro extensions.

See Also

[REVOKE](#), [ALTER DEFAULT PRIVILEGES](#)

INSERT

INSERT — create new rows in a table

Synopsis

```
[ WITH [ RECURSIVE ] with_query [, ...] ]
INSERT INTO table_name [ AS alias ] [ ( column_name [, ...] ) ]
    [ OVERRIDING { SYSTEM | USER } VALUE ]
    { DEFAULT VALUES | VALUES ( { expression | DEFAULT } [, ...] ) [, ...] | query }
    [ ON CONFLICT [ conflict_target ] conflict_action ]
    [ RETURNING { * | output_expression [ [ AS ] output_name ] } [, ...] ]
```

where *conflict_target* can be one of:

```
( { index_column_name | ( index_expression ) } [ COLLATE collation ] [ opclass ]
[, ...] ) [ WHERE index_predicate ]
ON CONSTRAINT constraint_name
```

and *conflict_action* is one of:

```
DO NOTHING
DO UPDATE SET { column_name = { expression | DEFAULT } |
                ( column_name [, ...] ) = [ ROW ] ( { expression | DEFAULT }
[, ...] ) |
                ( column_name [, ...] ) = ( sub-SELECT )
                } [, ...]
[ WHERE condition ]
```

Description

INSERT inserts new rows into a table. One can insert one or more rows specified by value expressions, or zero or more rows resulting from a query.

The target column names can be listed in any order. If no list of column names is given at all, the default is all the columns of the table in their declared order; or the first *N* column names, if there are only *N* columns supplied by the VALUES clause or *query*. The values supplied by the VALUES clause or *query* are associated with the explicit or implicit column list left-to-right.

Each column not present in the explicit or implicit column list will be filled with a default value, either its declared default value or null if there is none.

If the expression for any column is not of the correct data type, automatic type conversion will be attempted.

INSERT into tables that lack unique indexes will not be blocked by concurrent activity. Tables with unique indexes might block if concurrent sessions perform actions that lock or modify rows matching the unique index values being inserted; the details are covered in [Section 63.5](#). ON CONFLICT can be used to specify an alternative action to raising a unique constraint or exclusion constraint violation error. (See [ON CONFLICT Clause](#) below.)

The optional RETURNING clause causes INSERT to compute and return value(s) based on each row actually inserted (or updated, if an ON CONFLICT DO UPDATE clause was used). This is primarily useful for obtaining values that were supplied by defaults, such as a serial sequence number. However, any expression using the table's columns is allowed. The syntax of the RETURNING list is identical to that of the output list of SELECT. Only rows that were successfully inserted or updated will be returned. For example, if a row was locked but not updated because an ON CONFLICT DO UPDATE ... WHERE clause *condition* was not satisfied, the row will not be returned.

You must have `INSERT` privilege on a table in order to insert into it. If `ON CONFLICT DO UPDATE` is present, `UPDATE` privilege on the table is also required.

If a column list is specified, you only need `INSERT` privilege on the listed columns. Similarly, when `ON CONFLICT DO UPDATE` is specified, you only need `UPDATE` privilege on the column(s) that are listed to be updated. However, `ON CONFLICT DO UPDATE` also requires `SELECT` privilege on any column whose values are read in the `ON CONFLICT DO UPDATE` expressions or *condition*.

Use of the `RETURNING` clause requires `SELECT` privilege on all columns mentioned in `RETURNING`. If you use the *query* clause to insert rows from a query, you of course need to have `SELECT` privilege on any table or column used in the query.

Parameters

Inserting

This section covers parameters that may be used when only inserting new rows. Parameters *exclusively* used with the `ON CONFLICT` clause are described separately.

with_query

The `WITH` clause allows you to specify one or more subqueries that can be referenced by name in the `INSERT` query. See [Section 7.8](#) and [SELECT](#) for details.

It is possible for the *query* (`SELECT` statement) to also contain a `WITH` clause. In such a case both sets of *with_query* can be referenced within the *query*, but the second one takes precedence since it is more closely nested.

table_name

The name (optionally schema-qualified) of an existing table.

alias

A substitute name for *table_name*. When an alias is provided, it completely hides the actual name of the table. This is particularly useful when `ON CONFLICT DO UPDATE` targets a table named `excluded`, since that will otherwise be taken as the name of the special table representing the row proposed for insertion.

column_name

The name of a column in the table named by *table_name*. The column name can be qualified with a subfield name or array subscript, if needed. (Inserting into only some fields of a composite column leaves the other fields null.) When referencing a column with `ON CONFLICT DO UPDATE`, do not include the table's name in the specification of a target column. For example, `INSERT INTO table_name ... ON CONFLICT DO UPDATE SET table_name.col = 1` is invalid (this follows the general behavior for `UPDATE`).

`OVERRIDING SYSTEM VALUE`

If this clause is specified, then any values supplied for identity columns will override the default sequence-generated values.

For an identity column defined as `GENERATED ALWAYS`, it is an error to insert an explicit value (other than `DEFAULT`) without specifying either `OVERRIDING SYSTEM VALUE` or `OVERRIDING USER VALUE`. (For an identity column defined as `GENERATED BY DEFAULT`, `OVERRIDING SYSTEM VALUE` is the normal behavior and specifying it does nothing, but Postgres Pro allows it as an extension.)

`OVERRIDING USER VALUE`

If this clause is specified, then any values supplied for identity columns are ignored and the default sequence-generated values are applied.

This clause is useful for example when copying values between tables. Writing `INSERT INTO tbl2 OVERRIDING USER VALUE SELECT * FROM tbl1` will copy from `tbl1` all columns that are not identity columns in `tbl2` while values for the identity columns in `tbl2` will be generated by the sequences associated with `tbl2`.

DEFAULT VALUES

All columns will be filled with their default values, as if `DEFAULT` were explicitly specified for each column. (An `OVERRIDING` clause is not permitted in this form.)

expression

An expression or value to assign to the corresponding column.

DEFAULT

The corresponding column will be filled with its default value. An identity column will be filled with a new value generated by the associated sequence. For a generated column, specifying this is permitted but merely specifies the normal behavior of computing the column from its generation expression.

query

A query (`SELECT` statement) that supplies the rows to be inserted. Refer to the [SELECT](#) statement for a description of the syntax.

output_expression

An expression to be computed and returned by the `INSERT` command after each row is inserted or updated. The expression can use any column names of the table named by *table_name*. Write `*` to return all columns of the inserted or updated row(s).

output_name

A name to use for a returned column.

ON CONFLICT **Clause**

The optional `ON CONFLICT` clause specifies an alternative action to raising a unique violation or exclusion constraint violation error. For each individual row proposed for insertion, either the insertion proceeds, or, if an *arbiter* constraint or index specified by *conflict_target* is violated, the alternative *conflict_action* is taken. `ON CONFLICT DO NOTHING` simply avoids inserting a row as its alternative action. `ON CONFLICT DO UPDATE` updates the existing row that conflicts with the row proposed for insertion as its alternative action.

conflict_target can perform *unique index inference*. When performing inference, it consists of one or more *index_column_name* columns and/or *index_expression* expressions, and an optional *index_predicate*. All *table_name* unique indexes that, without regard to order, contain exactly the *conflict_target*-specified columns/expressions are inferred (chosen) as arbiter indexes. If an *index_predicate* is specified, it must, as a further requirement for inference, satisfy arbiter indexes. Note that this means a non-partial unique index (a unique index without a predicate) will be inferred (and thus used by `ON CONFLICT`) if such an index satisfying every other criteria is available. If an attempt at inference is unsuccessful, an error is raised.

`ON CONFLICT DO UPDATE` guarantees an atomic `INSERT` or `UPDATE` outcome; provided there is no independent error, one of those two outcomes is guaranteed, even under high concurrency. This is also known as *UPSERT* — “UPDATE or INSERT”.

conflict_target

Specifies which conflicts `ON CONFLICT` takes the alternative action on by choosing *arbiter indexes*. Either performs *unique index inference*, or names a constraint explicitly. For `ON CONFLICT DO NOTHING`, it is optional to specify a *conflict_target*; when omitted, conflicts with all usable constraints (and unique indexes) are handled. For `ON CONFLICT DO UPDATE`, a *conflict_target* *must* be provided.

conflict_action

conflict_action specifies an alternative `ON CONFLICT` action. It can be either `DO NOTHING`, or a `DO UPDATE` clause specifying the exact details of the `UPDATE` action to be performed in case of a conflict. The `SET` and `WHERE` clauses in `ON CONFLICT DO UPDATE` have access to the existing row using the table's name (or an alias), and to the row proposed for insertion using the special `excluded` table. `SELECT` privilege is required on any column in the target table where corresponding `excluded` columns are read.

Note that the effects of all per-row `BEFORE INSERT` triggers are reflected in `excluded` values, since those effects may have contributed to the row being excluded from insertion.

index_column_name

The name of a *table_name* column. Used to infer arbiter indexes. Follows `CREATE INDEX` format. `SELECT` privilege on *index_column_name* is required.

index_expression

Similar to *index_column_name*, but used to infer expressions on *table_name* columns appearing within index definitions (not simple columns). Follows `CREATE INDEX` format. `SELECT` privilege on any column appearing within *index_expression* is required.

collation

When specified, mandates that corresponding *index_column_name* or *index_expression* use a particular collation in order to be matched during inference. Typically this is omitted, as collations usually do not affect whether or not a constraint violation occurs. Follows `CREATE INDEX` format.

opclass

When specified, mandates that corresponding *index_column_name* or *index_expression* use particular operator class in order to be matched during inference. Typically this is omitted, as the *equality* semantics are often equivalent across a type's operator classes anyway, or because it's sufficient to trust that the defined unique indexes have the pertinent definition of equality. Follows `CREATE INDEX` format.

index_predicate

Used to allow inference of partial unique indexes. Any indexes that satisfy the predicate (which need not actually be partial indexes) can be inferred. Follows `CREATE INDEX` format. `SELECT` privilege on any column appearing within *index_predicate* is required.

constraint_name

Explicitly specifies an arbiter *constraint* by name, rather than inferring a constraint or index.

condition

An expression that returns a value of type `boolean`. Only rows for which this expression returns `true` will be updated, although all rows will be locked when the `ON CONFLICT DO UPDATE` action is taken. Note that *condition* is evaluated last, after a conflict has been identified as a candidate to update.

Note that exclusion constraints are not supported as arbiters with `ON CONFLICT DO UPDATE`. In all cases, only `NOT DEFERRABLE` constraints and unique indexes are supported as arbiters.

`INSERT` with an `ON CONFLICT DO UPDATE` clause is a “deterministic” statement. This means that the command will not be allowed to affect any single existing row more than once; a cardinality violation error will be raised when this situation arises. Rows proposed for insertion should not duplicate each other in terms of attributes constrained by an arbiter index or constraint.

Note that it is currently not supported for the `ON CONFLICT DO UPDATE` clause of an `INSERT` applied to a partitioned table to update the partition key of a conflicting row such that it requires the row be moved to a new partition.

Tip

It is often preferable to use unique index inference rather than naming a constraint directly using `ON CONFLICT ON CONSTRAINT constraint_name`. Inference will continue to work correctly when the underlying index is replaced by another more or less equivalent index in an overlapping way, for example when using `CREATE UNIQUE INDEX ... CONCURRENTLY` before dropping the index being replaced.

Outputs

On successful completion, an `INSERT` command returns a command tag of the form

```
INSERT oid count
```

The *count* is the number of rows inserted or updated. *oid* is always 0 (it used to be the OID assigned to the inserted row if *count* was exactly one and the target table was declared `WITH OIDS` and 0 otherwise, but creating a table `WITH OIDS` is not supported anymore).

If the `INSERT` command contains a `RETURNING` clause, the result will be similar to that of a `SELECT` statement containing the columns and values defined in the `RETURNING` list, computed over the row(s) inserted or updated by the command.

Notes

If the specified table is a partitioned table, each row is routed to the appropriate partition and inserted into it. If the specified table is a partition, an error will occur if one of the input rows violates the partition constraint.

You may also wish to consider using `MERGE`, since that allows mixing `INSERT`, `UPDATE`, and `DELETE` within a single statement. See [MERGE](#).

Examples

Insert a single row into table `films`:

```
INSERT INTO films VALUES
('UA502', 'Bananas', 105, '1971-07-13', 'Comedy', '82 minutes');
```

In this example, the `len` column is omitted and therefore it will have the default value:

```
INSERT INTO films (code, title, did, date_prod, kind)
VALUES ('T_601', 'Yojimbo', 106, '1961-06-16', 'Drama');
```

This example uses the `DEFAULT` clause for the date columns rather than specifying a value:

```
INSERT INTO films VALUES
('UA502', 'Bananas', 105, DEFAULT, 'Comedy', '82 minutes');
INSERT INTO films (code, title, did, date_prod, kind)
VALUES ('T_601', 'Yojimbo', 106, DEFAULT, 'Drama');
```

To insert a row consisting entirely of default values:

```
INSERT INTO films DEFAULT VALUES;
```

To insert multiple rows using the multirow `VALUES` syntax:

```
INSERT INTO films (code, title, did, date_prod, kind) VALUES
('B6717', 'Tampopo', 110, '1985-02-10', 'Comedy'),
('HG120', 'The Dinner Game', 140, DEFAULT, 'Comedy');
```

This example inserts some rows into table `films` from a table `tmp_films` with the same column layout as `films`:

```
INSERT INTO films SELECT * FROM tmp_films WHERE date_prod < '2004-05-07';
```

This example inserts into array columns:

```
-- Create an empty 3x3 gameboard for noughts-and-crosses
INSERT INTO tictactoe (game, board[1:3][1:3])
    VALUES (1, '{{" ", " ", " "},{ " ", " ", " "},{ " ", " ", " "}}');
-- The subscripts in the above example aren't really needed
INSERT INTO tictactoe (game, board)
    VALUES (2, '{{X, " ", " "},{ " ", O, " "},{ " ", X, " "}}');
```

Insert a single row into table distributors, returning the sequence number generated by the DEFAULT clause:

```
INSERT INTO distributors (did, dname) VALUES (DEFAULT, 'XYZ Widgets')
    RETURNING did;
```

Increment the sales count of the salesperson who manages the account for Acme Corporation, and record the whole updated row along with current time in a log table:

```
WITH upd AS (
    UPDATE employees SET sales_count = sales_count + 1 WHERE id =
        (SELECT sales_person FROM accounts WHERE name = 'Acme Corporation')
    RETURNING *
)
INSERT INTO employees_log SELECT *, current_timestamp FROM upd;
```

Insert or update new distributors as appropriate. Assumes a unique index has been defined that constrains values appearing in the did column. Note that the special excluded table is used to reference values originally proposed for insertion:

```
INSERT INTO distributors (did, dname)
    VALUES (5, 'Gizmo Transglobal'), (6, 'Associated Computing, Inc')
    ON CONFLICT (did) DO UPDATE SET dname = EXCLUDED.dname;
```

Insert a distributor, or do nothing for rows proposed for insertion when an existing, excluded row (a row with a matching constrained column or columns after before row insert triggers fire) exists. Example assumes a unique index has been defined that constrains values appearing in the did column:

```
INSERT INTO distributors (did, dname) VALUES (7, 'Redline GmbH')
    ON CONFLICT (did) DO NOTHING;
```

Insert or update new distributors as appropriate. Example assumes a unique index has been defined that constrains values appearing in the did column. WHERE clause is used to limit the rows actually updated (any existing row not updated will still be locked, though):

```
-- Don't update existing distributors based in a certain ZIP code
INSERT INTO distributors AS d (did, dname) VALUES (8, 'Anvil Distribution')
    ON CONFLICT (did) DO UPDATE
    SET dname = EXCLUDED.dname || ' (formerly ' || d.dname || ' )'
    WHERE d.zipcode <> '21201';

-- Name a constraint directly in the statement (uses associated
-- index to arbitrate taking the DO NOTHING action)
INSERT INTO distributors (did, dname) VALUES (9, 'Antwerp Design')
    ON CONFLICT ON CONSTRAINT distributors_pkey DO NOTHING;
```

Insert new distributor if possible; otherwise DO NOTHING. Example assumes a unique index has been defined that constrains values appearing in the did column on a subset of rows where the is_active Boolean column evaluates to true:

```
-- This statement could infer a partial unique index on "did"
-- with a predicate of "WHERE is_active", but it could also
```

```
-- just use a regular unique constraint on "did"
INSERT INTO distributors (did, dname) VALUES (10, 'Conrad International')
    ON CONFLICT (did) WHERE is_active DO NOTHING;
```

Compatibility

INSERT conforms to the SQL standard, except that the RETURNING clause is a Postgres Pro extension, as is the ability to use WITH with INSERT, and the ability to specify an alternative action with ON CONFLICT. Also, the case in which a column name list is omitted, but not all the columns are filled from the VALUES clause or *query*, is disallowed by the standard. If you prefer a more SQL standard conforming statement than ON CONFLICT, see [MERGE](#).

The SQL standard specifies that OVERRIDING SYSTEM VALUE can only be specified if an identity column that is generated always exists. Postgres Pro allows the clause in any case and ignores it if it is not applicable.

Possible limitations of the *query* clause are documented under [SELECT](#).

LISTEN

LISTEN — listen for a notification

Synopsis

```
LISTEN channel
```

Description

LISTEN registers the current session as a listener on the notification channel named *channel*. If the current session is already registered as a listener for this notification channel, nothing is done.

Whenever the command NOTIFY *channel* is invoked, either by this session or another one connected to the same database, all the sessions currently listening on that notification channel are notified, and each will in turn notify its connected client application.

A session can be unregistered for a given notification channel with the UNLISTEN command. A session's listen registrations are automatically cleared when the session ends.

The method a client application must use to detect notification events depends on which Postgres Pro application programming interface it uses. With the libpq library, the application issues LISTEN as an ordinary SQL command, and then must periodically call the function PQnotifies to find out whether any notification events have been received. Other interfaces such as libpgtcl provide higher-level methods for handling notify events; indeed, with libpgtcl the application programmer should not even issue LISTEN or UNLISTEN directly. See the documentation for the interface you are using for more details.

Parameters

channel

Name of a notification channel (any identifier).

Notes

LISTEN takes effect at transaction commit. If LISTEN or UNLISTEN is executed within a transaction that later rolls back, the set of notification channels being listened to is unchanged.

A transaction that has executed LISTEN cannot be prepared for two-phase commit.

There is a race condition when first setting up a listening session: if concurrently-committing transactions are sending notify events, exactly which of those will the newly listening session receive? The answer is that the session will receive all events committed after an instant during the transaction's commit step. But that is slightly later than any database state that the transaction could have observed in queries. This leads to the following rule for using LISTEN: first execute (and commit!) that command, then in a new transaction inspect the database state as needed by the application logic, then rely on notifications to find out about subsequent changes to the database state. The first few received notifications might refer to updates already observed in the initial database inspection, but this is usually harmless.

NOTIFY contains a more extensive discussion of the use of LISTEN and NOTIFY.

Examples

Configure and execute a listen/notify sequence from psql:

```
LISTEN virtual;
NOTIFY virtual;
Asynchronous notification "virtual" received from server process with PID 8448.
```

Compatibility

There is no `LISTEN` statement in the SQL standard.

See Also

[NOTIFY](#), [UNLISTEN](#), [max_notify_queue_pages](#)

LOAD

LOAD — load a shared library file

Synopsis

```
LOAD 'filename'
```

Description

This command loads a shared library file into the Postgres Pro server's address space. If the file has been loaded already, the command does nothing. Shared library files that contain C functions are automatically loaded whenever one of their functions is called. Therefore, an explicit `LOAD` is usually only needed to load a library that modifies the server's behavior through “hooks” rather than providing a set of functions.

The library file name is typically given as just a bare file name, which is sought in the server's library search path (set by [dynamic_library_path](#)). Alternatively it can be given as a full path name. In either case the platform's standard shared library file name extension may be omitted. See [Section 39.10.1](#) for more information on this topic.

Non-superusers can only apply `LOAD` to library files located in `$libdir/plugins/` — the specified *filename* must begin with exactly that string. (It is the database administrator's responsibility to ensure that only “safe” libraries are installed there.)

Compatibility

`LOAD` is a Postgres Pro extension.

See Also

[CREATE FUNCTION](#)

LOCK

LOCK — lock a table

Synopsis

For distributed systems (sharded and global tables):

```
LOCK [ TABLE ] [ ONLY ] name [ * ] [, ...] [ IN lockmode MODE ] [ NOWAIT ]
```

where *lockmode* is one of:

```
ACCESS SHARE | SHARE UPDATE EXCLUSIVE | SHARE | EXCLUSIVE | ACCESS EXCLUSIVE
```

For other system types:

```
LOCK [ TABLE ] [ ONLY ] name [ * ] [, ...] [ IN lockmode MODE ] [ NOWAIT ]
```

where *lockmode* is one of:

```
ACCESS SHARE | ROW SHARE | ROW EXCLUSIVE | SHARE UPDATE EXCLUSIVE  
| SHARE | SHARE ROW EXCLUSIVE | EXCLUSIVE | ACCESS EXCLUSIVE
```

Description

`LOCK TABLE` obtains a table-level lock, waiting if necessary for any conflicting locks to be released. If `NOWAIT` is specified, `LOCK TABLE` does not wait to acquire the desired lock: if it cannot be acquired immediately, the command is aborted and an error is emitted. Once obtained, the lock is held for the remainder of the current transaction. (There is no `UNLOCK TABLE` command; locks are always released at transaction end.) For global and sharded tables the lock is on a table level (or one part of the table) for all cluster nodes.

Note that you cannot `LOCK` local and sharded or global tables simultaneously.

When a view is locked, all relations appearing in the view definition query are also locked recursively with the same lock mode.

When acquiring locks automatically for commands that reference tables, Postgres Pro always uses the least restrictive lock mode possible. `LOCK TABLE` provides for cases when you might need more restrictive locking. For example, suppose an application runs a transaction at the `READ COMMITTED` isolation level and needs to ensure that data in a table remains stable for the duration of the transaction. To achieve this you could obtain `SHARE` lock mode over the table before querying. This will prevent concurrent data changes and ensure subsequent reads of the table see a stable view of committed data, because `SHARE` lock mode conflicts with the `ROW EXCLUSIVE` lock acquired by writers, and your `LOCK TABLE name IN SHARE MODE` statement will wait until any concurrent holders of `ROW EXCLUSIVE` mode locks commit or roll back. Thus, once you obtain the lock, there are no uncommitted writes outstanding; furthermore none can begin until you release the lock.

To achieve a similar effect when running a transaction at the `REPEATABLE READ` or `SERIALIZABLE` isolation level, you have to execute the `LOCK TABLE` statement before executing any `SELECT` or data modification statement. A `REPEATABLE READ` or `SERIALIZABLE` transaction's view of data will be frozen when its first `SELECT` or data modification statement begins. A `LOCK TABLE` later in the transaction will still prevent concurrent writes — but it won't ensure that what the transaction reads corresponds to the latest committed values.

If a transaction of this sort is going to change the data in the table, then it should use `SHARE ROW EXCLUSIVE` lock mode instead of `SHARE` mode. This ensures that only one transaction of this type runs at a time. Without this, a deadlock is possible: two transactions might both acquire `SHARE` mode, and then be unable to also acquire `ROW EXCLUSIVE` mode to actually perform their updates. (Note that a transaction's

own locks never conflict, so a transaction can acquire `ROW EXCLUSIVE` mode when it holds `SHARE` mode — but not if anyone else holds `SHARE` mode.) To avoid deadlocks, make sure all transactions acquire locks on the same objects in the same order, and if multiple lock modes are involved for a single object, then transactions should always acquire the most restrictive mode first.

More information about the lock modes and locking strategies can be found in [Section 13.3](#).

Parameters

name

The name (optionally schema-qualified) of an existing table to lock. If `ONLY` is specified before the table name, only that table is locked. If `ONLY` is not specified, the table and all its descendant tables (if any) are locked. Optionally, `*` can be specified after the table name to explicitly indicate that descendant tables are included.

The command `LOCK TABLE a, b;` is equivalent to `LOCK TABLE a; LOCK TABLE b;`. The tables are locked one-by-one in the order specified in the `LOCK TABLE` command.

lockmode

The lock mode specifies which locks this lock conflicts with. Lock modes are described in [Section 13.3](#).

If no lock mode is specified, then `ACCESS EXCLUSIVE`, the most restrictive mode, is used.

`NOWAIT`

Specifies that `LOCK TABLE` should not wait for any conflicting locks to be released: if the specified lock(s) cannot be acquired immediately without waiting, the transaction is aborted.

Notes

To lock a table, the user must have the right privilege for the specified *lockmode*. If the user has `MAINTAIN`, `UPDATE`, `DELETE`, or `TRUNCATE` privileges on the table, any *lockmode* is permitted. If the user has `INSERT` privileges on the table, `ROW EXCLUSIVE MODE` (or a less-conflicting mode as described in [Section 13.3](#)) is permitted. If a user has `SELECT` privileges on the table, `ACCESS SHARE MODE` is permitted.

The user performing the lock on the view must have the corresponding privilege on the view. In addition, by default, the view's owner must have the relevant privileges on the underlying base relations, whereas the user performing the lock does not need any permissions on the underlying base relations. However, if the view has `security_invoker` set to `true` (see [CREATE VIEW](#)), the user performing the lock, rather than the view owner, must have the relevant privileges on the underlying base relations.

`LOCK TABLE` is useless outside a transaction block: the lock would remain held only to the completion of the statement. Therefore Postgres Pro reports an error if `LOCK` is used outside a transaction block. Use [BEGIN](#) and [COMMIT](#) (or [ROLLBACK](#)) to define a transaction block.

`LOCK TABLE` only deals with table-level locks, and so the mode names involving `ROW` are all misnomers. These mode names should generally be read as indicating the intention of the user to acquire row-level locks within the locked table. Also, `ROW EXCLUSIVE mode` is a shareable table lock. Keep in mind that all the lock modes have identical semantics so far as `LOCK TABLE` is concerned, differing only in the rules about which modes conflict with which. For information on how to acquire an actual row-level lock, see [Section 13.3.2](#) and [The Locking Clause](#) in the [SELECT](#) documentation.

Examples

Obtain a `SHARE` lock on a primary key table when going to perform inserts into a foreign key table:

```
BEGIN WORK;  
LOCK TABLE films IN SHARE MODE;  
SELECT id FROM films
```

```
WHERE name = 'Star Wars: Episode I - The Phantom Menace';
-- Do ROLLBACK if record was not returned
INSERT INTO films_user_comments VALUES
  (_id_, 'GREAT! I was waiting for it for so long!');
COMMIT WORK;
```

Take a `SHARE ROW EXCLUSIVE` lock on a primary key table when going to perform a delete operation:

```
BEGIN WORK;
LOCK TABLE films IN SHARE ROW EXCLUSIVE MODE;
DELETE FROM films_user_comments WHERE id IN
  (SELECT id FROM films WHERE rating < 5);
DELETE FROM films WHERE rating < 5;
COMMIT WORK;
```

Compatibility

There is no `LOCK TABLE` in the SQL standard, which instead uses `SET TRANSACTION` to specify concurrency levels on transactions. Postgres Pro supports that too; see [SET TRANSACTION](#) for details.

Except for `ACCESS SHARE`, `ACCESS EXCLUSIVE`, and `SHARE UPDATE EXCLUSIVE` lock modes, the Postgres Pro lock modes and the `LOCK TABLE` syntax are compatible with those present in Oracle.

MERGE

MERGE — conditionally insert, update, or delete rows of a table

Synopsis

```
[ WITH with_query [, ...] ]
MERGE INTO [ ONLY ] target_table_name [ * ] [ [ AS ] target_alias ]
USING data_source ON join_condition
when_clause [...]
[ RETURNING { * | output_expression [ [ AS ] output_name ] } [, ...] ]

where data_source is:

{ [ ONLY ] source_table_name [ * ] | ( source_query ) } [ [ AS ] source_alias ]

and when_clause is:

{ WHEN MATCHED [ AND condition ] THEN { merge_update | merge_delete | DO NOTHING } |
  WHEN NOT MATCHED BY SOURCE [ AND condition ] THEN { merge_update | merge_delete | DO
  NOTHING } |
  WHEN NOT MATCHED [ BY TARGET ] [ AND condition ] THEN { merge_insert | DO NOTHING } }

and merge_insert is:

INSERT [( column_name [, ...] )]
[ OVERRIDING { SYSTEM | USER } VALUE ]
{ VALUES ( { expression | DEFAULT } [, ...] ) | DEFAULT VALUES }

and merge_update is:

UPDATE SET { column_name = { expression | DEFAULT } |
            ( column_name [, ...] ) = [ ROW ] ( { expression | DEFAULT } [, ...] ) |
            ( column_name [, ...] ) = ( sub-SELECT )
            } [, ...]

and merge_delete is:

DELETE
```

Description

MERGE performs actions that modify rows in the target table identified as *target_table_name*, using the *data_source*. MERGE provides a single SQL statement that can conditionally INSERT, UPDATE or DELETE rows, a task that would otherwise require multiple procedural language statements.

First, the MERGE command performs a join from *data_source* to the target table producing zero or more candidate change rows. For each candidate change row, the status of MATCHED, NOT MATCHED BY SOURCE, or NOT MATCHED [BY TARGET] is set just once, after which WHEN clauses are evaluated in the order specified. For each candidate change row, the first clause to evaluate as true is executed. No more than one WHEN clause is executed for any candidate change row.

MERGE actions have the same effect as regular UPDATE, INSERT, or DELETE commands of the same names. The syntax of those commands is different, notably that there is no WHERE clause and no table name is specified. All actions refer to the target table, though modifications to other tables may be made using triggers.

When `DO NOTHING` is specified, the source row is skipped. Since actions are evaluated in their specified order, `DO NOTHING` can be handy to skip non-interesting source rows before more fine-grained handling.

The optional `RETURNING` clause causes `MERGE` to compute and return value(s) based on each row inserted, updated, or deleted. Any expression using the source or target table's columns, or the `merge_action()` function can be computed. When an `INSERT` or `UPDATE` action is performed, the new values of the target table's columns are used. When a `DELETE` is performed, the old values of the target table's columns are used. The syntax of the `RETURNING` list is identical to that of the output list of `SELECT`.

There is no separate `MERGE` privilege. If you specify an update action, you must have the `UPDATE` privilege on the column(s) of the target table that are referred to in the `SET` clause. If you specify an insert action, you must have the `INSERT` privilege on the target table. If you specify a delete action, you must have the `DELETE` privilege on the target table. If you specify a `DO NOTHING` action, you must have the `SELECT` privilege on at least one column of the target table. You will also need `SELECT` privilege on any column(s) of the `data_source` and of the target table referred to in any `condition` (including `join_condition`) or `expression`. Privileges are tested once at statement start and are checked whether or not particular `WHEN` clauses are executed.

`MERGE` is not supported if the target table is a materialized view, foreign table, or if it has any rules defined on it.

Parameters

with_query

The `WITH` clause allows you to specify one or more subqueries that can be referenced by name in the `MERGE` query. See [Section 7.8](#) and [SELECT](#) for details. Note that `WITH RECURSIVE` is not supported by `MERGE`.

target_table_name

The name (optionally schema-qualified) of the target table or view to merge into. If `ONLY` is specified before a table name, matching rows are updated or deleted in the named table only. If `ONLY` is not specified, matching rows are also updated or deleted in any tables inheriting from the named table. Optionally, `*` can be specified after the table name to explicitly indicate that descendant tables are included. The `ONLY` keyword and `*` option do not affect insert actions, which always insert into the named table only.

If *target_table_name* is a view, it must either be automatically updatable with no `INSTEAD OF` triggers, or it must have `INSTEAD OF` triggers for every type of action (`INSERT`, `UPDATE`, and `DELETE`) specified in the `WHEN` clauses. Views with rules are not supported.

target_alias

A substitute name for the target table. When an alias is provided, it completely hides the actual name of the table. For example, given `MERGE INTO foo AS f`, the remainder of the `MERGE` statement must refer to this table as `f` not `foo`.

source_table_name

The name (optionally schema-qualified) of the source table, view, or transition table. If `ONLY` is specified before the table name, matching rows are included from the named table only. If `ONLY` is not specified, matching rows are also included from any tables inheriting from the named table. Optionally, `*` can be specified after the table name to explicitly indicate that descendant tables are included.

source_query

A query (`SELECT` statement or `VALUES` statement) that supplies the rows to be merged into the target table. Refer to the [SELECT](#) statement or [VALUES](#) statement for a description of the syntax.

source_alias

A substitute name for the data source. When an alias is provided, it completely hides the actual name of the table or the fact that a query was issued.

join_condition

join_condition is an expression resulting in a value of type `boolean` (similar to a `WHERE` clause) that specifies which rows in the *data_source* match rows in the target table.

Warning

Only columns from the target table that attempt to match *data_source* rows should appear in *join_condition*. *join_condition* subexpressions that only reference the target table's columns can affect which action is taken, often in surprising ways.

If both `WHEN NOT MATCHED BY SOURCE` and `WHEN NOT MATCHED [BY TARGET]` clauses are specified, the `MERGE` command will perform a `FULL` join between *data_source* and the target table. For this to work, at least one *join_condition* subexpression must use an operator that can support a hash join, or all of the subexpressions must use operators that can support a merge join.

when_clause

At least one `WHEN` clause is required.

The `WHEN` clause may specify `WHEN MATCHED`, `WHEN NOT MATCHED BY SOURCE`, or `WHEN NOT MATCHED [BY TARGET]`. Note that the SQL standard only defines `WHEN MATCHED` and `WHEN NOT MATCHED` (which is defined to mean no matching target row). `WHEN NOT MATCHED BY SOURCE` is an extension to the SQL standard, as is the option to append `BY TARGET` to `WHEN NOT MATCHED`, to make its meaning more explicit.

If the `WHEN` clause specifies `WHEN MATCHED` and the candidate change row matches a row in the *data_source* to a row in the target table, the `WHEN` clause is executed if the *condition* is absent or it evaluates to `true`.

If the `WHEN` clause specifies `WHEN NOT MATCHED BY SOURCE` and the candidate change row represents a row in the target table that does not match a row in the *data_source*, the `WHEN` clause is executed if the *condition* is absent or it evaluates to `true`.

If the `WHEN` clause specifies `WHEN NOT MATCHED [BY TARGET]` and the candidate change row represents a row in the *data_source* that does not match a row in the target table, the `WHEN` clause is executed if the *condition* is absent or it evaluates to `true`.

condition

An expression that returns a value of type `boolean`. If this expression for a `WHEN` clause returns `true`, then the action for that clause is executed for that row.

A condition on a `WHEN MATCHED` clause can refer to columns in both the source and the target relations. A condition on a `WHEN NOT MATCHED BY SOURCE` clause can only refer to columns from the target relation, since by definition there is no matching source row. A condition on a `WHEN NOT MATCHED [BY TARGET]` clause can only refer to columns from the source relation, since by definition there is no matching target row. Only the system attributes from the target table are accessible.

merge_insert

The specification of an `INSERT` action that inserts one row into the target table. The target column names can be listed in any order. If no list of column names is given at all, the default is all the columns of the table in their declared order.

Each column not present in the explicit or implicit column list will be filled with a default value, either its declared default value or null if there is none.

If the target table is a partitioned table, each row is routed to the appropriate partition and inserted into it. If the target table is a partition, an error will occur if any input row violates the partition constraint.

Column names may not be specified more than once. `INSERT` actions cannot contain sub-selects.

Only one `VALUES` clause can be specified. The `VALUES` clause can only refer to columns from the source relation, since by definition there is no matching target row.

merge_update

The specification of an `UPDATE` action that updates the current row of the target table. Column names may not be specified more than once.

Neither a table name nor a `WHERE` clause are allowed.

merge_delete

Specifies a `DELETE` action that deletes the current row of the target table. Do not include the table name or any other clauses, as you would normally do with a `DELETE` command.

column_name

The name of a column in the target table. The column name can be qualified with a subfield name or array subscript, if needed. (Inserting into only some fields of a composite column leaves the other fields null.) Do not include the table's name in the specification of a target column.

OVERRIDING SYSTEM VALUE

Without this clause, it is an error to specify an explicit value (other than `DEFAULT`) for an identity column defined as `GENERATED ALWAYS`. This clause overrides that restriction.

OVERRIDING USER VALUE

If this clause is specified, then any values supplied for identity columns defined as `GENERATED BY DEFAULT` are ignored and the default sequence-generated values are applied.

DEFAULT VALUES

All columns will be filled with their default values. (An `OVERRIDING` clause is not permitted in this form.)

expression

An expression to assign to the column. If used in a `WHEN MATCHED` clause, the expression can use values from the original row in the target table, and values from the `data_source` row. If used in a `WHEN NOT MATCHED BY SOURCE` clause, the expression can only use values from the original row in the target table. If used in a `WHEN NOT MATCHED [BY TARGET]` clause, the expression can only use values from the `data_source` row.

DEFAULT

Set the column to its default value (which will be `NULL` if no specific default expression has been assigned to it).

sub-SELECT

A `SELECT` sub-query that produces as many output columns as are listed in the parenthesized column list preceding it. The sub-query must yield no more than one row when executed. If it yields one row, its column values are assigned to the target columns; if it yields no rows, `NULL` values are assigned to the target columns. If used in a `WHEN MATCHED` clause, the sub-query can refer to values from the

original row in the target table, and values from the `data_source` row. If used in a `WHEN NOT MATCHED BY SOURCE` clause, the sub-query can only refer to values from the original row in the target table.

output_expression

An expression to be computed and returned by the `MERGE` command after each row is changed (whether inserted, updated, or deleted). The expression can use any columns of the source or target tables, or the `merge_action()` function to return additional information about the action executed.

Writing `*` will return all columns from the source table, followed by all columns from the target table. Often this will lead to a lot of duplication, since it is common for the source and target tables to have a lot of the same columns. This can be avoided by qualifying the `*` with the name or alias of the source or target table.

output_name

A name to use for a returned column.

Outputs

On successful completion, a `MERGE` command returns a command tag of the form

`MERGE total_count`

The `total_count` is the total number of rows changed (whether inserted, updated, or deleted). If `total_count` is 0, no rows were changed in any way.

If the `MERGE` command contains a `RETURNING` clause, the result will be similar to that of a `SELECT` statement containing the columns and values defined in the `RETURNING` list, computed over the row(s) inserted, updated, or deleted by the command.

Notes

The following steps take place during the execution of `MERGE`.

1. Perform any `BEFORE STATEMENT` triggers for all actions specified, whether or not their `WHEN` clauses match.
2. Perform a join from source to target table. The resulting query will be optimized normally and will produce a set of candidate change rows. For each candidate change row,
 - a. Evaluate whether each row is `MATCHED`, `NOT MATCHED BY SOURCE`, or `NOT MATCHED [BY TARGET]`.
 - b. Test each `WHEN` condition in the order specified until one returns true.
 - c. When a condition returns true, perform the following actions:
 - i. Perform any `BEFORE ROW` triggers that fire for the action's event type.
 - ii. Perform the specified action, invoking any check constraints on the target table.
 - iii. Perform any `AFTER ROW` triggers that fire for the action's event type.If the target relation is a view with `INSTEAD OF ROW` triggers for the action's event type, they are used to perform the action instead.
3. Perform any `AFTER STATEMENT` triggers for actions specified, whether or not they actually occur. This is similar to the behavior of an `UPDATE` statement that modifies no rows.

In summary, statement triggers for an event type (say, `INSERT`) will be fired whenever we *specify* an action of that kind. In contrast, row-level triggers will fire only for the specific event type being *executed*. So a `MERGE` command might fire statement triggers for both `UPDATE` and `INSERT`, even though only `UPDATE` row triggers were fired.

You should ensure that the join produces at most one candidate change row for each target row. In other words, a target row shouldn't join to more than one data source row. If it does, then only one of the candidate change rows will be used to modify the target row; later attempts to modify the row will

cause an error. This can also occur if row triggers make changes to the target table and the rows so modified are then subsequently also modified by `MERGE`. If the repeated action is an `INSERT`, this will cause a uniqueness violation, while a repeated `UPDATE` or `DELETE` will cause a cardinality violation; the latter behavior is required by the SQL standard. This differs from historical PostgreSQL behavior of joins in `UPDATE` and `DELETE` statements where second and subsequent attempts to modify the same row are simply ignored.

If a `WHEN` clause omits an `AND` sub-clause, it becomes the final reachable clause of that kind (`MATCHED`, `NOT MATCHED BY SOURCE`, or `NOT MATCHED [BY TARGET]`). If a later `WHEN` clause of that kind is specified it would be provably unreachable and an error is raised. If no final reachable clause is specified of either kind, it is possible that no action will be taken for a candidate change row.

The order in which rows are generated from the data source is indeterminate by default. A *source_query* can be used to specify a consistent ordering, if required, which might be needed to avoid deadlocks between concurrent transactions.

When `MERGE` is run concurrently with other commands that modify the target table, the usual transaction isolation rules apply; see [Section 13.2](#) for an explanation on the behavior at each isolation level. You may also wish to consider using `INSERT ... ON CONFLICT` as an alternative statement which offers the ability to run an `UPDATE` if a concurrent `INSERT` occurs. There are a variety of differences and restrictions between the two statement types and they are not interchangeable.

Examples

Perform maintenance on `customer_accounts` based upon new `recent_transactions`.

```
MERGE INTO customer_account ca
USING recent_transactions t
ON t.customer_id = ca.customer_id
WHEN MATCHED THEN
    UPDATE SET balance = balance + transaction_value
WHEN NOT MATCHED THEN
    INSERT (customer_id, balance)
    VALUES (t.customer_id, t.transaction_value);
```

Notice that this would be exactly equivalent to the following statement because the `MATCHED` result does not change during execution.

```
MERGE INTO customer_account ca
USING (SELECT customer_id, transaction_value FROM recent_transactions) AS t
ON t.customer_id = ca.customer_id
WHEN MATCHED THEN
    UPDATE SET balance = balance + transaction_value
WHEN NOT MATCHED THEN
    INSERT (customer_id, balance)
    VALUES (t.customer_id, t.transaction_value);
```

Attempt to insert a new stock item along with the quantity of stock. If the item already exists, instead update the stock count of the existing item. Don't allow entries that have zero stock. Return details of all changes made.

```
MERGE INTO wines w
USING wine_stock_changes s
ON s.wine_name = w.wine_name
WHEN NOT MATCHED AND s.stock_delta > 0 THEN
    INSERT VALUES(s.wine_name, s.stock_delta)
WHEN MATCHED AND w.stock + s.stock_delta > 0 THEN
    UPDATE SET stock = w.stock + s.stock_delta
WHEN MATCHED THEN
    DELETE
```

```
RETURNING merge_action(), w.*;
```

The `wine_stock_changes` table might be, for example, a temporary table recently loaded into the database.

Update `wines` based on a replacement wine list, inserting rows for any new stock, updating modified stock entries, and deleting any wines not present in the new list.

```
MERGE INTO wines w
USING new_wine_list s
ON s.wine_name = w.wine_name
WHEN NOT MATCHED BY TARGET THEN
    INSERT VALUES(s.wine_name, s.stock)
WHEN MATCHED AND w.stock != s.stock THEN
    UPDATE SET stock = s.stock
WHEN NOT MATCHED BY SOURCE THEN
    DELETE;
```

Compatibility

This command conforms to the SQL standard.

The `WITH` clause, `BY SOURCE` and `BY TARGET` qualifiers to `WHEN NOT MATCHED, DO NOTHING` action, and `RETURNING` clause are extensions to the SQL standard.

MOVE

MOVE — position a cursor

Synopsis

```
MOVE [ direction ] [ FROM | IN ] cursor_name
```

where *direction* can be one of:

```
NEXT
PRIOR
FIRST
LAST
ABSOLUTE count
RELATIVE count
count
ALL
FORWARD
FORWARD count
FORWARD ALL
BACKWARD
BACKWARD count
BACKWARD ALL
```

Description

MOVE repositions a cursor without retrieving any data. MOVE works exactly like the FETCH command, except it only positions the cursor and does not return rows.

The parameters for the MOVE command are identical to those of the FETCH command; refer to [FETCH](#) for details on syntax and usage.

Outputs

On successful completion, a MOVE command returns a command tag of the form

```
MOVE count
```

The *count* is the number of rows that a FETCH command with the same parameters would have returned (possibly zero).

Examples

```
BEGIN WORK;
DECLARE liahona CURSOR FOR SELECT * FROM films;

-- Skip the first 5 rows:
MOVE FORWARD 5 IN liahona;
MOVE 5

-- Fetch the 6th row from the cursor liahona:
FETCH 1 FROM liahona;
  code | title  | did | date_prod | kind  | len
-----+-----+-----+-----+-----+-----
  P_303 | 48 Hrs | 103 | 1982-10-22 | Action | 01:37
(1 row)
```

```
-- Close the cursor liahona and end the transaction:  
CLOSE liahona;  
COMMIT WORK;
```

Compatibility

There is no `MOVE` statement in the SQL standard.

See Also

[CLOSE](#), [DECLARE](#), [FETCH](#)

NOTIFY

NOTIFY — generate a notification

Synopsis

```
NOTIFY channel [ , payload ]
```

Description

The `NOTIFY` command sends a notification event together with an optional “payload” string to each client application that has previously executed `LISTEN channel` for the specified channel name in the current database. Notifications are visible to all users.

`NOTIFY` provides a simple interprocess communication mechanism for a collection of processes accessing the same Postgres Pro database. A payload string can be sent along with the notification, and higher-level mechanisms for passing structured data can be built by using tables in the database to pass additional data from notifier to listener(s).

The information passed to the client for a notification event includes the notification channel name, the notifying session's server process PID, and the payload string, which is an empty string if it has not been specified.

It is up to the database designer to define the channel names that will be used in a given database and what each one means. Commonly, the channel name is the same as the name of some table in the database, and the notify event essentially means, “I changed this table, take a look at it to see what's new”. But no such association is enforced by the `NOTIFY` and `LISTEN` commands. For example, a database designer could use several different channel names to signal different sorts of changes to a single table. Alternatively, the payload string could be used to differentiate various cases.

When `NOTIFY` is used to signal the occurrence of changes to a particular table, a useful programming technique is to put the `NOTIFY` in a statement trigger that is triggered by table updates. In this way, notification happens automatically when the table is changed, and the application programmer cannot accidentally forget to do it.

`NOTIFY` interacts with SQL transactions in some important ways. Firstly, if a `NOTIFY` is executed inside a transaction, the notify events are not delivered until and unless the transaction is committed. This is appropriate, since if the transaction is aborted, all the commands within it have had no effect, including `NOTIFY`. But it can be disconcerting if one is expecting the notification events to be delivered immediately. Secondly, if a listening session receives a notification signal while it is within a transaction, the notification event will not be delivered to its connected client until just after the transaction is completed (either committed or aborted). Again, the reasoning is that if a notification were delivered within a transaction that was later aborted, one would want the notification to be undone somehow — but the server cannot “take back” a notification once it has sent it to the client. So notification events are only delivered between transactions. The upshot of this is that applications using `NOTIFY` for real-time signaling should try to keep their transactions short.

If the same channel name is signaled multiple times with identical payload strings within the same transaction, only one instance of the notification event is delivered to listeners. On the other hand, notifications with distinct payload strings will always be delivered as distinct notifications. Similarly, notifications from different transactions will never get folded into one notification. Except for dropping later instances of duplicate notifications, `NOTIFY` guarantees that notifications from the same transaction get delivered in the order they were sent. It is also guaranteed that messages from different transactions are delivered in the order in which the transactions committed.

It is common for a client that executes `NOTIFY` to be listening on the same notification channel itself. In that case it will get back a notification event, just like all the other listening sessions. Depending on the application logic, this could result in useless work, for example, reading a database table to find the

same updates that that session just wrote out. It is possible to avoid such extra work by noticing whether the notifying session's server process PID (supplied in the notification event message) is the same as one's own session's PID (available from `libpq`). When they are the same, the notification event is one's own work bouncing back, and can be ignored.

Parameters

channel

Name of the notification channel to be signaled (any identifier).

payload

The “payload” string to be communicated along with the notification. This must be specified as a simple string literal. In the default configuration it must be shorter than 8000 bytes. (If binary data or large amounts of information need to be communicated, it's best to put it in a database table and send the key of the record.)

Notes

There is a queue that holds notifications that have been sent but not yet processed by all listening sessions. If this queue becomes full, transactions calling `NOTIFY` will fail at commit. The queue is quite large (8GB in a standard installation) and should be sufficiently sized for almost every use case. However, no cleanup can take place if a session executes `LISTEN` and then enters a transaction for a very long time. Once the queue is half full you will see warnings in the log file pointing you to the session that is preventing cleanup. In this case you should make sure that this session ends its current transaction so that cleanup can proceed.

The function `pg_notification_queue_usage` returns the fraction of the queue that is currently occupied by pending notifications. See [Section 9.27](#) for more information.

A transaction that has executed `NOTIFY` cannot be prepared for two-phase commit.

pg_notify

To send a notification you can also use the function `pg_notify(text, text)`. The function takes the channel name as the first argument and the payload as the second. The function is much easier to use than the `NOTIFY` command if you need to work with non-constant channel names and payloads.

Examples

Configure and execute a listen/notify sequence from `psql`:

```
LISTEN virtual;
NOTIFY virtual;
Asynchronous notification "virtual" received from server process with PID 8448.
NOTIFY virtual, 'This is the payload';
Asynchronous notification "virtual" with payload "This is the payload" received from
server process with PID 8448.

LISTEN foo;
SELECT pg_notify('fo' || 'o', 'pay' || 'load');
Asynchronous notification "foo" with payload "payload" received from server process
with PID 14728.
```

Compatibility

There is no `NOTIFY` statement in the SQL standard.

See Also

[LISTEN](#), [UNLISTEN](#), [max_notify_queue_pages](#)

PREPARE

PREPARE — prepare a statement for execution

Synopsis

```
PREPARE name [ ( data_type [, ...] ) ] AS statement
```

Description

PREPARE creates a prepared statement. A prepared statement is a server-side object that can be used to optimize performance. When the PREPARE statement is executed, the specified statement is parsed, analyzed, and rewritten. When an EXECUTE command is subsequently issued, the prepared statement is planned and executed. This division of labor avoids repetitive parse analysis work, while allowing the execution plan to depend on the specific parameter values supplied.

Prepared statements can take parameters: values that are substituted into the statement when it is executed. When creating the prepared statement, refer to parameters by position, using \$1, \$2, etc. A corresponding list of parameter data types can optionally be specified. When a parameter's data type is not specified or is declared as `unknown`, the type is inferred from the context in which the parameter is first referenced (if possible). When executing the statement, specify the actual values for these parameters in the EXECUTE statement. Refer to [EXECUTE](#) for more information about that.

Prepared statements only last for the duration of the current database session. When the session ends, the prepared statement is forgotten, so it must be recreated before being used again. This also means that a single prepared statement cannot be used by multiple simultaneous database clients; however, each client can create their own prepared statement to use. Prepared statements can be manually cleaned up using the [DEALLOCATE](#) command.

Prepared statements potentially have the largest performance advantage when a single session is being used to execute a large number of similar statements. The performance difference will be particularly significant if the statements are complex to plan or rewrite, e.g., if the query involves a join of many tables or requires the application of several rules. If the statement is relatively simple to plan and rewrite but relatively expensive to execute, the performance advantage of prepared statements will be less noticeable.

Parameters

name

An arbitrary name given to this particular prepared statement. It must be unique within a single session and is subsequently used to execute or deallocate a previously prepared statement.

data_type

The data type of a parameter to the prepared statement. If the data type of a particular parameter is unspecified or is specified as `unknown`, it will be inferred from the context in which the parameter is first referenced. To refer to the parameters in the prepared statement itself, use \$1, \$2, etc.

statement

Any SELECT, INSERT, UPDATE, DELETE, MERGE, or VALUES statement.

Notes

A prepared statement can be executed with either a *generic plan* or a *custom plan*. A generic plan is the same across all executions, while a custom plan is generated for a specific execution using the parameter values given in that call. Use of a generic plan avoids planning overhead, but in some situations a custom plan will be much more efficient to execute because the planner can make use of knowledge of the

parameter values. (Of course, if the prepared statement has no parameters, then this is moot and a generic plan is always used.)

By default (that is, when [plan_cache_mode](#) is set to `auto`), the server will automatically choose whether to use a generic or custom plan for a prepared statement that has parameters. The current rule for this is that the first five executions are done with custom plans and the average estimated cost of those plans is calculated. Then a generic plan is created and its estimated cost is compared to the average custom-plan cost. Subsequent executions use the generic plan if its cost is not so much higher than the average custom-plan cost as to make repeated replanning seem preferable.

This heuristic can be overridden, forcing the server to use either generic or custom plans, by setting `plan_cache_mode` to `force_generic_plan` or `force_custom_plan` respectively. This setting is primarily useful if the generic plan's cost estimate is badly off for some reason, allowing it to be chosen even though its actual cost is much more than that of a custom plan.

By default, generic plans for only 64 most recent prepared statements are kept in memory; older statements have to be parsed and analyzed again if they are called later. This limit is enforced by the [plan_cache_lru_size](#) configuration parameter, which enables you to control the amount of memory used by prepared statements. If you would like to keep plans in memory for all prepared statements as long as possible, set this parameter to 0.

To examine the query plan Postgres Pro is using for a prepared statement, use [EXPLAIN](#), for example

```
EXPLAIN EXECUTE name(parameter_values);
```

If a generic plan is in use, it will contain parameter symbols $\$n$, while a custom plan will have the supplied parameter values substituted into it.

For more information on query planning and the statistics collected by Postgres Pro for that purpose, see the [ANALYZE](#) documentation.

Although the main point of a prepared statement is to avoid repeated parse analysis and planning of the statement, Postgres Pro will force re-analysis and re-planning of the statement before using it whenever database objects used in the statement have undergone definitional (DDL) changes or their planner statistics have been updated since the previous use of the prepared statement. Also, if the value of [search_path](#) changes from one use to the next, the statement will be re-parsed using the new `search_path`. (This latter behavior is new as of PostgreSQL 9.3.) These rules make use of a prepared statement semantically almost equivalent to re-submitting the same query text over and over, but with a performance benefit if no object definitions are changed, especially if the best plan remains the same across uses. An example of a case where the semantic equivalence is not perfect is that if the statement refers to a table by an unqualified name, and then a new table of the same name is created in a schema appearing earlier in the `search_path`, no automatic re-parse will occur since no object used in the statement changed. However, if some other change forces a re-parse, the new table will be referenced in subsequent uses.

You can see all prepared statements available in the session by querying the [pg_prepared_statements](#) system view.

Examples

Create a prepared statement for an `INSERT` statement, and then execute it:

```
PREPARE fooplan (int, text, bool, numeric) AS
  INSERT INTO foo VALUES($1, $2, $3, $4);
EXECUTE fooplan(1, 'Hunter Valley', 't', 200.00);
```

Create a prepared statement for a `SELECT` statement, and then execute it:

```
PREPARE usrrptplan (int) AS
  SELECT * FROM users u, logs l WHERE u.usrid=$1 AND u.usrid=l.usrid
  AND l.date = $2;
```

```
EXECUTE usrrptplan(1, current_date);
```

In this example, the data type of the second parameter is not specified, so it is inferred from the context in which \$2 is used.

Compatibility

The SQL standard includes a `PREPARE` statement, but it is only for use in embedded SQL. This version of the `PREPARE` statement also uses a somewhat different syntax.

See Also

[DEALLOCATE](#), [EXECUTE](#)

PREPARE TRANSACTION

PREPARE TRANSACTION — prepare the current transaction for two-phase commit

Synopsis

```
PREPARE TRANSACTION transaction_id
```

Description

PREPARE TRANSACTION prepares the current transaction for two-phase commit. After this command, the transaction is no longer associated with the current session; instead, its state is fully stored on disk, and there is a very high probability that it can be committed successfully, even if a database crash occurs before the commit is requested.

Once prepared, a transaction can later be committed or rolled back with `COMMIT PREPARED` or `ROLLBACK PREPARED`, respectively. Those commands can be issued from any session, not only the one that executed the original transaction.

From the point of view of the issuing session, PREPARE TRANSACTION is not unlike a ROLLBACK command: after executing it, there is no active current transaction, and the effects of the prepared transaction are no longer visible. (The effects will become visible again if the transaction is committed.)

If the PREPARE TRANSACTION command fails for any reason, it becomes a ROLLBACK: the current transaction is canceled.

Parameters

transaction_id

An arbitrary identifier that later identifies this transaction for `COMMIT PREPARED` or `ROLLBACK PREPARED`. The identifier must be written as a string literal, and must be less than 200 bytes long. It must not be the same as the identifier used for any currently prepared transaction.

Notes

PREPARE TRANSACTION is not intended for use in applications or interactive sessions. Its purpose is to allow an external transaction manager to perform atomic global transactions across multiple databases or other transactional resources. Unless you're writing a transaction manager, you probably shouldn't be using PREPARE TRANSACTION.

This command must be used inside a transaction block. Use `BEGIN` to start one.

It is not currently allowed to PREPARE a transaction that has executed any operations involving temporary tables or the session's temporary namespace, created any cursors `WITH HOLD`, or executed `LISTEN`, `UNLISTEN`, or `NOTIFY`. Those features are too tightly tied to the current session to be useful in a transaction to be prepared.

If the transaction modified any run-time parameters with `SET` (without the `LOCAL` option), those effects persist after PREPARE TRANSACTION, and will not be affected by any later `COMMIT PREPARED` or `ROLLBACK PREPARED`. Thus, in this one respect PREPARE TRANSACTION acts more like `COMMIT` than `ROLLBACK`.

All currently available prepared transactions are listed in the `pg_prepared_xacts` system view.

Caution

It is unwise to leave transactions in the prepared state for a long time. This will interfere with the ability of `VACUUM` to reclaim storage. Keep in mind also that the transaction continues to hold

whatever locks it held. The intended usage of the feature is that a prepared transaction will normally be committed or rolled back as soon as an external transaction manager has verified that other databases are also prepared to commit.

If you have not set up an external transaction manager to track prepared transactions and ensure they get closed out promptly, it is best to keep the prepared-transaction feature disabled by setting [max_prepared_transactions](#) to zero. This will prevent accidental creation of prepared transactions that might then be forgotten and eventually cause problems.

Examples

Prepare the current transaction for two-phase commit, using `foobar` as the transaction identifier:

```
PREPARE TRANSACTION 'foobar';
```

Compatibility

`PREPARE TRANSACTION` is a Postgres Pro extension. It is intended for use by external transaction management systems, some of which are covered by standards (such as X/Open XA), but the SQL side of those systems is not standardized.

See Also

[COMMIT PREPARED](#), [ROLLBACK PREPARED](#)

REASSIGN OWNED

REASSIGN OWNED — change the ownership of database objects owned by a database role

Synopsis

```
REASSIGN OWNED BY { old_role | CURRENT_ROLE | CURRENT_USER | SESSION_USER } [, ...]
                TO { new_role | CURRENT_ROLE | CURRENT_USER | SESSION_USER }
```

Description

REASSIGN OWNED instructs the system to change the ownership of database objects owned by any of the *old_roles* to *new_role*.

Parameters

old_role

The name of a role. The ownership of all the objects within the current database, and of all shared objects (databases, tablespaces), owned by this role will be reassigned to *new_role*.

new_role

The name of the role that will be made the new owner of the affected objects.

Notes

REASSIGN OWNED is often used to prepare for the removal of one or more roles. Because REASSIGN OWNED does not affect objects within other databases, it is usually necessary to execute this command in each database that contains objects owned by a role that is to be removed.

REASSIGN OWNED requires membership on both the source role(s) and the target role.

The [DROP OWNED](#) command is an alternative that simply drops all the database objects owned by one or more roles.

The REASSIGN OWNED command does not affect any privileges granted to the *old_roles* on objects that are not owned by them. Likewise, it does not affect default privileges created with ALTER DEFAULT PRIVILEGES. Use DROP OWNED to revoke such privileges.

See [Section 21.4](#) for more discussion.

Compatibility

The REASSIGN OWNED command is a Postgres Pro extension.

See Also

[DROP OWNED](#), [DROP ROLE](#), [ALTER DATABASE](#)

REFRESH MATERIALIZED VIEW

REFRESH MATERIALIZED VIEW — replace the contents of a materialized view

Synopsis

```
REFRESH MATERIALIZED VIEW [ CONCURRENTLY ] name
[ WITH [ NO ] DATA ]
```

Description

REFRESH MATERIALIZED VIEW completely replaces the contents of a materialized view. To execute this command you must have the `MAINTAIN` privilege on the materialized view. The old contents are discarded. If `WITH DATA` is specified (or defaults) the backing query is executed to provide the new data, and the materialized view is left in a scannable state. If `WITH NO DATA` is specified no new data is generated and the materialized view is left in an unscannable state.

`CONCURRENTLY` and `WITH NO DATA` may not be specified together.

Parameters

`CONCURRENTLY`

Refresh the materialized view without locking out concurrent selects on the materialized view. Without this option a refresh which affects a lot of rows will tend to use fewer resources and complete more quickly, but could block other connections which are trying to read from the materialized view. This option may be faster in cases where a small number of rows are affected.

This option is only allowed if there is at least one `UNIQUE` index on the materialized view which uses only column names and includes all rows; that is, it must not be an expression index or include a `WHERE` clause.

This option can only be used when the materialized view is already populated.

Even with this option only one `REFRESH` at a time may run against any one materialized view.

name

The name (optionally schema-qualified) of the materialized view to refresh.

Notes

If there is an `ORDER BY` clause in the materialized view's defining query, the original contents of the materialized view will be ordered that way; but `REFRESH MATERIALIZED VIEW` does not guarantee to preserve that ordering.

While `REFRESH MATERIALIZED VIEW` is running, the [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`.

Examples

This command will replace the contents of the materialized view called `order_summary` using the query from the materialized view's definition, and leave it in a scannable state:

```
REFRESH MATERIALIZED VIEW order_summary;
```

This command will free storage associated with the materialized view `annual_statistics_basis` and leave it in an unscannable state:

```
REFRESH MATERIALIZED VIEW annual_statistics_basis WITH NO DATA;
```


Compatibility

REFRESH MATERIALIZED VIEW is a Postgres Pro extension.

See Also

[CREATE MATERIALIZED VIEW](#), [ALTER MATERIALIZED VIEW](#), [DROP MATERIALIZED VIEW](#)

REINDEX

REINDEX — rebuild indexes

Synopsis

```
REINDEX [ ( option [, ...] ) ] { INDEX | TABLE | SCHEMA } [ CONCURRENTLY ] name
REINDEX [ ( option [, ...] ) ] { DATABASE | SYSTEM } [ CONCURRENTLY ] [ name ]
```

where *option* can be one of:

```
CONCURRENTLY [ boolean ]
TABLESPACE new_tablespace
VERBOSE [ boolean ]
```

Description

REINDEX rebuilds an index using the data stored in the index's table, replacing the old copy of the index. There are several scenarios in which to use **REINDEX**:

- An index has become corrupted, and no longer contains valid data. Although in theory this should never happen, in practice indexes can become corrupted due to software bugs or hardware failures. **REINDEX** provides a recovery method.
- An index has become “bloated”, that is it contains many empty or nearly-empty pages. This can occur with B-tree indexes in Postgres Pro under certain uncommon access patterns. **REINDEX** provides a way to reduce the space consumption of the index by writing a new version of the index without the dead pages. See [Section 24.3](#) for more information.
- You have altered a storage parameter (such as `fillfactor`) for an index, and wish to ensure that the change has taken full effect.
- If an index build fails with the `CONCURRENTLY` option, this index is left as “invalid”. Such indexes are useless but it can be convenient to use **REINDEX** to rebuild them. Note that only **REINDEX INDEX** is able to perform a concurrent build on an invalid index.

Parameters

INDEX

Recreate the specified index. This form of **REINDEX** cannot be executed inside a transaction block when used with a partitioned index.

TABLE

Recreate all indexes of the specified table. If the table has a secondary “TOAST” table, that is reindexed as well. This form of **REINDEX** cannot be executed inside a transaction block when used with a partitioned table.

SCHEMA

Recreate all indexes of the specified schema. If a table of this schema has a secondary “TOAST” table, that is reindexed as well. Indexes on shared system catalogs are also processed. This form of **REINDEX** cannot be executed inside a transaction block.

DATABASE

Recreate all indexes within the current database, except system catalogs. Indexes on system catalogs are not processed. This form of **REINDEX** cannot be executed inside a transaction block.

SYSTEM

Recreate all indexes on system catalogs within the current database. Indexes on shared system catalogs are included. Indexes on user tables are not processed. This form of `REINDEX` cannot be executed inside a transaction block.

name

The name of the specific index, table, or database to be reindexed. Index and table names can be schema-qualified. Presently, `REINDEX DATABASE` and `REINDEX SYSTEM` can only reindex the current database. Their parameter is optional, and it must match the current database's name.

CONCURRENTLY

When this option is used, Postgres Pro will rebuild the index without taking any locks that prevent concurrent inserts, updates, or deletes on the table; whereas a standard index rebuild locks out writes (but not reads) on the table until it's done. There are several caveats to be aware of when using this option — see [Rebuilding Indexes Concurrently](#) below.

For temporary tables, `REINDEX` is always non-concurrent, as no other session can access them, and non-concurrent reindex is cheaper.

TABLESPACE

Specifies that indexes will be rebuilt on a new tablespace.

VERBOSE

Prints a progress report as each index is reindexed.

boolean

Specifies whether the selected option should be turned on or off. You can write `TRUE`, `ON`, or `1` to enable the option, and `FALSE`, `OFF`, or `0` to disable it. The *boolean* value can also be omitted, in which case `TRUE` is assumed.

new_tablespace

The tablespace where indexes will be rebuilt.

Notes

If you suspect corruption of an index on a user table, you can simply rebuild that index, or all indexes on the table, using `REINDEX INDEX` or `REINDEX TABLE`.

Things are more difficult if you need to recover from corruption of an index on a system table. In this case it's important for the system to not have used any of the suspect indexes itself. (Indeed, in this sort of scenario you might find that server processes are crashing immediately at start-up, due to reliance on the corrupted indexes.) To recover safely, the server must be started with the `-P` option, which prevents it from using indexes for system catalog lookups.

One way to do this is to shut down the server and start a single-user Postgres Pro server with the `-P` option included on its command line. Then, `REINDEX DATABASE`, `REINDEX SYSTEM`, `REINDEX TABLE`, or `REINDEX INDEX` can be issued, depending on how much you want to reconstruct. If in doubt, use `REINDEX SYSTEM` to select reconstruction of all system indexes in the database. Then quit the single-user server session and restart the regular server. See the [postgres](#) reference page for more information about how to interact with the single-user server interface.

Alternatively, a regular server session can be started with `-P` included in its command line options. The method for doing this varies across clients, but in all libpq-based clients, it is possible to set the `PGOPTIONS` environment variable to `-P` before starting the client. Note that while this method does not

require locking out other clients, it might still be wise to prevent other users from connecting to the damaged database until repairs have been completed.

`REINDEX` is similar to a drop and recreate of the index in that the index contents are rebuilt from scratch. However, the locking considerations are rather different. `REINDEX` locks out writes but not reads of the index's parent table. It also takes an `ACCESS EXCLUSIVE` lock on the specific index being processed, which will block reads that attempt to use that index. In particular, the query planner tries to take an `ACCESS SHARE` lock on every index of the table, regardless of the query, and so `REINDEX` blocks virtually any queries except for some prepared queries whose plan has been cached and which don't use this very index. In contrast, `DROP INDEX` momentarily takes an `ACCESS EXCLUSIVE` lock on the parent table, blocking both writes and reads. The subsequent `CREATE INDEX` locks out writes but not reads; since the index is not there, no read will attempt to use it, meaning that there will be no blocking but reads might be forced into expensive sequential scans.

While `REINDEX` is running, the [search_path](#) is temporarily changed to `pg_catalog, pg_temp`.

Reindexing a single index or table requires having the `MAINTAIN` privilege on the table. Note that while `REINDEX` on a partitioned index or table requires having the `MAINTAIN` privilege on the partitioned table, such commands skip the privilege checks when processing the individual partitions. Reindexing a schema or database requires being the owner of that schema or database or having privileges of the [pg_maintain](#) role. Note specifically that it's thus possible for non-superusers to rebuild indexes of tables owned by other users. However, as a special exception, `REINDEX DATABASE`, `REINDEX SCHEMA`, and `REINDEX SYSTEM` will skip indexes on shared catalogs unless the user has the `MAINTAIN` privilege on the catalog.

Reindexing partitioned indexes or partitioned tables is supported with `REINDEX INDEX` or `REINDEX TABLE`, respectively. Each partition of the specified partitioned relation is reindexed in a separate transaction. Those commands cannot be used inside a transaction block when working on a partitioned table or index.

When using the `TABLESPACE` clause with `REINDEX` on a partitioned index or table, only the tablespace references of the leaf partitions are updated. As partitioned indexes are not updated, it is recommended to separately use `ALTER TABLE ONLY` on them so as any new partitions attached inherit the new tablespace. On failure, it may not have moved all the indexes to the new tablespace. Re-running the command will rebuild all the leaf partitions and move previously-unprocessed indexes to the new tablespace.

If `SCHEMA`, `DATABASE` or `SYSTEM` is used with `TABLESPACE`, system relations are skipped and a single `WARNING` will be generated. Indexes on TOAST tables are rebuilt, but not moved to the new tablespace.

Rebuilding Indexes Concurrently

Rebuilding an index can interfere with regular operation of a database. Normally Postgres Pro locks the table whose index is rebuilt against writes and performs the entire index build with a single scan of the table. Other transactions can still read the table, but if they try to insert, update, or delete rows in the table they will block until the index rebuild is finished. This could have a severe effect if the system is a live production database. Very large tables can take many hours to be indexed, and even for smaller tables, an index rebuild can lock out writers for periods that are unacceptably long for a production system.

Postgres Pro supports rebuilding indexes with minimum locking of writes. This method is invoked by specifying the `CONCURRENTLY` option of `REINDEX`. When this option is used, Postgres Pro must perform two scans of the table for each index that needs to be rebuilt and wait for termination of all existing transactions that could potentially use the index. This method requires more total work than a standard index rebuild and takes significantly longer to complete as it needs to wait for unfinished transactions that might modify the index. However, since it allows normal operations to continue while the index is being rebuilt, this method is useful for rebuilding indexes in a production environment. Of course, the extra CPU, memory and I/O load imposed by the index rebuild may slow down other operations.

The following steps occur in a concurrent reindex. Each step is run in a separate transaction. If there are multiple indexes to be rebuilt, then each step loops through all the indexes before moving to the next step.

1. A new transient index definition is added to the catalog `pg_index`. This definition will be used to replace the old index. A `SHARE UPDATE EXCLUSIVE` lock at session level is taken on the indexes being reindexed as well as their associated tables to prevent any schema modification while processing.
2. A first pass to build the index is done for each new index. Once the index is built, its flag `pg_index.indisready` is switched to “true” to make it ready for inserts, making it visible to other sessions once the transaction that performed the build is finished. This step is done in a separate transaction for each index.
3. Then a second pass is performed to add tuples that were added while the first pass was running. This step is also done in a separate transaction for each index.
4. All the constraints that refer to the index are changed to refer to the new index definition, and the names of the indexes are changed. At this point, `pg_index.indisvalid` is switched to “true” for the new index and to “false” for the old, and a cache invalidation is done causing all sessions that referenced the old index to be invalidated.
5. The old indexes have `pg_index.indisready` switched to “false” to prevent any new tuple insertions, after waiting for running queries that might reference the old index to complete.
6. The old indexes are dropped. The `SHARE UPDATE EXCLUSIVE` session locks for the indexes and the table are released.

If a problem arises while rebuilding the indexes, such as a uniqueness violation in a unique index, the `REINDEX` command will fail but leave behind an “invalid” new index in addition to the pre-existing one. This index will be ignored for querying purposes because it might be incomplete; however it will still consume update overhead. The `psql \d` command will report such an index as `INVALID`:

```
postgres=# \d tab
           Table "public.tab"
  Column | Type   | Modifiers
-----+-----+-----
 col    | integer |
Indexes:
    "idx" btree (col)
    "idx_ccnew" btree (col) INVALID
```

If the index marked `INVALID` is suffixed `ccnew`, then it corresponds to the transient index created during the concurrent operation, and the recommended recovery method is to drop it using `DROP INDEX`, then attempt `REINDEX CONCURRENTLY` again. If the invalid index is instead suffixed `ccold`, it corresponds to the original index which could not be dropped; the recommended recovery method is to just drop said index, since the rebuild proper has been successful.

Regular index builds permit other regular index builds on the same table to occur simultaneously, but only one concurrent index build can occur on a table at a time. In both cases, no other types of schema modification on the table are allowed meanwhile. Another difference is that a regular `REINDEX TABLE` or `REINDEX INDEX` command can be performed within a transaction block, but `REINDEX CONCURRENTLY` cannot.

Like any long-running transaction, `REINDEX` on a table can affect which tuples can be removed by concurrent `VACUUM` on any other table.

`REINDEX SYSTEM` does not support `CONCURRENTLY` since system catalogs cannot be reindexed concurrently.

Furthermore, indexes for exclusion constraints cannot be reindexed concurrently. If such an index is named directly in this command, an error is raised. If a table or database with exclusion constraint indexes is reindexed concurrently, those indexes will be skipped. (It is possible to reindex such indexes without the `CONCURRENTLY` option.)

Each backend running `REINDEX` will report its progress in the `pg_stat_progress_create_index` view. See [Section 27.4.4](#) for details.

Examples

Rebuild a single index:

```
REINDEX INDEX my_index;
```

Rebuild all the indexes on the table `my_table`:

```
REINDEX TABLE my_table;
```

Rebuild all indexes in a particular database, without trusting the system indexes to be valid already:

```
$ export PGOPTIONS="-P"
```

```
$ psql broken_db
```

```
...
```

```
broken_db=> REINDEX DATABASE broken_db;
```

```
broken_db=> \q
```

Rebuild indexes for a table, without blocking read and write operations on involved relations while reindexing is in progress:

```
REINDEX TABLE CONCURRENTLY my_broken_table;
```

Compatibility

There is no `REINDEX` command in the SQL standard.

See Also

[CREATE INDEX](#), [DROP INDEX](#), [reindexdb](#), [Section 27.4.4](#)

RELEASE SAVEPOINT

RELEASE SAVEPOINT — release a previously defined savepoint

Synopsis

```
RELEASE [ SAVEPOINT ] savepoint_name
```

Description

RELEASE SAVEPOINT releases the named savepoint and all active savepoints that were created after the named savepoint, and frees their resources. All changes made since the creation of the savepoint that didn't already get rolled back are merged into the transaction or savepoint that was active when the named savepoint was created. Changes made after RELEASE SAVEPOINT will also be part of this active transaction or savepoint.

Parameters

savepoint_name

The name of the savepoint to release.

Notes

Specifying a savepoint name that was not previously defined is an error.

It is not possible to release a savepoint when the transaction is in an aborted state; to do that, use [ROLLBACK TO SAVEPOINT](#).

If multiple savepoints have the same name, only the most recently defined unreleased one is released. Repeated commands will release progressively older savepoints.

Examples

To establish and later release a savepoint:

```
BEGIN;  
  INSERT INTO table1 VALUES (3);  
  SAVEPOINT my_savepoint;  
  INSERT INTO table1 VALUES (4);  
  RELEASE SAVEPOINT my_savepoint;  
COMMIT;
```

The above transaction will insert both 3 and 4.

A more complex example with multiple nested subtransactions:

```
BEGIN;  
  INSERT INTO table1 VALUES (1);  
  SAVEPOINT sp1;  
  INSERT INTO table1 VALUES (2);  
  SAVEPOINT sp2;  
  INSERT INTO table1 VALUES (3);  
  RELEASE SAVEPOINT sp2;  
  INSERT INTO table1 VALUES (4)); -- generates an error
```

In this example, the application requests the release of the savepoint `sp2`, which inserted 3. This changes the insert's transaction context to `sp1`. When the statement attempting to insert value 4 generates an error, the insertion of 2 and 4 are lost because they are in the same, now-rolled back savepoint, and value 3 is in the same transaction context. The application can now only choose one of these two commands, since all other commands will be ignored:

```
ROLLBACK;  
ROLLBACK TO SAVEPOINT sp1;
```

Choosing `ROLLBACK` will abort everything, including value 1, whereas `ROLLBACK TO SAVEPOINT sp1` will retain value 1 and allow the transaction to continue.

Compatibility

This command conforms to the SQL standard. The standard specifies that the key word `SAVEPOINT` is mandatory, but Postgres Pro allows it to be omitted.

See Also

[BEGIN](#), [COMMIT](#), [ROLLBACK](#), [ROLLBACK TO SAVEPOINT](#), [SAVEPOINT](#)

RESET

RESET — restore the value of a run-time parameter to the default value

Synopsis

```
RESET configuration_parameter
RESET ALL
```

Description

RESET restores run-time parameters to their default values. RESET is an alternative spelling for

```
SET configuration_parameter TO DEFAULT
```

Refer to [SET](#) for details.

The default value is defined as the value that the parameter would have had, if no SET had ever been issued for it in the current session. The actual source of this value might be a compiled-in default, the configuration file, command-line options, or per-database or per-user default settings. This is subtly different from defining it as “the value that the parameter had at session start”, because if the value came from the configuration file, it will be reset to whatever is specified by the configuration file now. See [Chapter 19](#) for details.

The transactional behavior of RESET is the same as SET: its effects will be undone by transaction rollback.

Parameters

configuration_parameter

Name of a settable run-time parameter. Available parameters are documented in [Chapter 19](#) and on the [SET](#) reference page.

ALL

Resets all settable run-time parameters to default values.

Examples

Set the `timezone` configuration variable to its default value:

```
RESET timezone;
```

Compatibility

RESET is a Postgres Pro extension.

See Also

[SET](#), [SHOW](#)

REVOKE

REVOKE — remove access privileges

Synopsis

```
REVOKE [ GRANT OPTION FOR ]
    { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER |
      MAINTAIN }
      [, ...] | ALL [ PRIVILEGES ] }
    ON { [ TABLE ] table_name [, ...]
        | ALL TABLES IN SCHEMA schema_name [, ...] }
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { { SELECT | INSERT | UPDATE | REFERENCES } ( column_name [, ...] )
      [, ...] | ALL [ PRIVILEGES ] ( column_name [, ...] ) }
    ON [ TABLE ] table_name [, ...]
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { { USAGE | SELECT | UPDATE }
      [, ...] | ALL [ PRIVILEGES ] }
    ON { SEQUENCE sequence_name [, ...]
        | ALL SEQUENCES IN SCHEMA schema_name [, ...] }
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { { CREATE | CONNECT | TEMPORARY | TEMP } [, ...] | ALL [ PRIVILEGES ] }
    ON DATABASE database_name [, ...]
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { USAGE | ALL [ PRIVILEGES ] }
    ON DOMAIN domain_name [, ...]
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { USAGE | ALL [ PRIVILEGES ] }
    ON FOREIGN DATA WRAPPER fdw_name [, ...]
    FROM role_specification [, ...]
    [ GRANTED BY role_specification ]
    [ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
    { USAGE | ALL [ PRIVILEGES ] }
    ON FOREIGN SERVER server_name [, ...]
```

```
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ EXECUTE | ALL [ PRIVILEGES ] }
ON { { FUNCTION | PROCEDURE | ROUTINE } function_name [ ( [ [ argmode ] [ arg_name
] arg_type [, ...] ] ) ] [, ...]
| ALL { FUNCTIONS | PROCEDURES | ROUTINES } IN SCHEMA schema_name [, ...] }
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ USAGE | ALL [ PRIVILEGES ] }
ON LANGUAGE lang_name [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ { SELECT | UPDATE } [, ...] | ALL [ PRIVILEGES ] }
ON LARGE OBJECT loid [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ { SET | ALTER SYSTEM } [, ...] | ALL [ PRIVILEGES ] }
ON PARAMETER configuration_parameter [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ { CREATE | USAGE } [, ...] | ALL [ PRIVILEGES ] }
ON SCHEMA schema_name [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ CREATE | ALL [ PRIVILEGES ] }
ON TABLESPACE tablespace_name [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ GRANT OPTION FOR ]
{ USAGE | ALL [ PRIVILEGES ] }
ON TYPE type_name [, ...]
FROM role_specification [, ...]
[ GRANTED BY role_specification ]
[ CASCADE | RESTRICT ]
```

```
REVOKE [ { ADMIN | INHERIT | SET } OPTION FOR ]
role_name [, ...] FROM role_specification [, ...]
[ GRANTED BY role_specification ]
```

```
[ CASCADE | RESTRICT ]
```

where *role_specification* can be:

```
[ GROUP ] role_name  
| PUBLIC  
| CURRENT_ROLE  
| CURRENT_USER  
| SESSION_USER
```

Description

The **REVOKE** command revokes previously granted privileges from one or more roles. The key word **PUBLIC** refers to the implicitly defined group of all roles.

See the description of the **GRANT** command for the meaning of the privilege types.

Note that any particular role will have the sum of privileges granted directly to it, privileges granted to any role it is presently a member of, and privileges granted to **PUBLIC**. Thus, for example, revoking **SELECT** privilege from **PUBLIC** does not necessarily mean that all roles have lost **SELECT** privilege on the object: those who have it granted directly or via another role will still have it. Similarly, revoking **SELECT** from a user might not prevent that user from using **SELECT** if **PUBLIC** or another membership role still has **SELECT** rights.

If **GRANT OPTION FOR** is specified, only the grant option for the privilege is revoked, not the privilege itself. Otherwise, both the privilege and the grant option are revoked.

If a user holds a privilege with grant option and has granted it to other users then the privileges held by those other users are called dependent privileges. If the privilege or the grant option held by the first user is being revoked and dependent privileges exist, those dependent privileges are also revoked if **CASCADE** is specified; if it is not, the revoke action will fail. This recursive revocation only affects privileges that were granted through a chain of users that is traceable to the user that is the subject of this **REVOKE** command. Thus, the affected users might effectively keep the privilege if it was also granted through other users.

When revoking privileges on a table, the corresponding column privileges (if any) are automatically revoked on each column of the table, as well. On the other hand, if a role has been granted privileges on a table, then revoking the same privileges from individual columns will have no effect.

When revoking membership in a role, **GRANT OPTION** is instead called **ADMIN OPTION**, but the behavior is similar. Note that, in releases prior to Postgres Pro 16, dependent privileges were not tracked for grants of role membership, and thus **CASCADE** had no effect for role membership. This is no longer the case. Note also that this form of the command does not allow the noise word **GROUP** in *role_specification*.

Just as **ADMIN OPTION** can be removed from an existing role grant, it is also possible to revoke **INHERIT OPTION** or **SET OPTION**. This is equivalent to setting the value of the corresponding option to **FALSE**.

Notes

A user can only revoke privileges that were granted directly by that user. If, for example, user A has granted a privilege with grant option to user B, and user B has in turn granted it to user C, then user A cannot revoke the privilege directly from C. Instead, user A could revoke the grant option from user B and use the **CASCADE** option so that the privilege is in turn revoked from user C. For another example, if both A and B have granted the same privilege to C, A can revoke their own grant but not B's grant, so C will still effectively have the privilege.

When a non-owner of an object attempts to **REVOKE** privileges on the object, the command will fail outright if the user has no privileges whatsoever on the object. As long as some privilege is available, the command will proceed, but it will revoke only those privileges for which the user has grant options. The

`REVOKE ALL PRIVILEGES` forms will issue a warning message if no grant options are held, while the other forms will issue a warning if grant options for any of the privileges specifically named in the command are not held. (In principle these statements apply to the object owner as well, but since the owner is always treated as holding all grant options, the cases can never occur.)

If a superuser chooses to issue a `GRANT` or `REVOKE` command, the command is performed as though it were issued by the owner of the affected object. (Since roles do not have owners, in the case of a `GRANT` of role membership, the command is performed as though it were issued by the bootstrap superuser.) Since all privileges ultimately come from the object owner (possibly indirectly via chains of grant options), it is possible for a superuser to revoke all privileges, but this might require use of `CASCADE` as stated above.

`REVOKE` can also be done by a role that is not the owner of the affected object, but is a member of the role that owns the object, or is a member of a role that holds privileges `WITH GRANT OPTION` on the object. In this case the command is performed as though it were issued by the containing role that actually owns the object or holds the privileges `WITH GRANT OPTION`. For example, if table `t1` is owned by role `g1`, of which role `u1` is a member, then `u1` can revoke privileges on `t1` that are recorded as being granted by `g1`. This would include grants made by `u1` as well as by other members of role `g1`.

If the role executing `REVOKE` holds privileges indirectly via more than one role membership path, it is unspecified which containing role will be used to perform the command. In such cases it is best practice to use `SET ROLE` to become the specific role you want to do the `REVOKE` as. Failure to do so might lead to revoking privileges other than the ones you intended, or not revoking anything at all.

See [Section 5.8](#) for more information about specific privilege types, as well as how to inspect objects' privileges.

Examples

Revoke insert privilege for the public on table `films`:

```
REVOKE INSERT ON films FROM PUBLIC;
```

Revoke all privileges from user `manuel` on view `kinds`:

```
REVOKE ALL PRIVILEGES ON kinds FROM manuel;
```

Note that this actually means “revoke all privileges that I granted”.

Revoke membership in role `admins` from user `joe`:

```
REVOKE admins FROM joe;
```

Compatibility

The compatibility notes of the [GRANT](#) command apply analogously to `REVOKE`. The keyword `RESTRICT` or `CASCADE` is required according to the standard, but Postgres Pro assumes `RESTRICT` by default.

See Also

[GRANT](#), [ALTER DEFAULT PRIVILEGES](#)

ROLLBACK

ROLLBACK — abort the current transaction

Synopsis

```
ROLLBACK [ AUTONOMOUS ] [ WORK | TRANSACTION ] [ AND [ NO ] CHAIN ]
```

Description

ROLLBACK rolls back the current transaction and causes all the updates made by the transaction to be discarded.

Parameters

AUTONOMOUS

Optional key word that can be used when aborting an autonomous transaction. For details on autonomous transactions, see [Chapter 16](#).

WORK

TRANSACTION

Optional key words. They have no effect.

AND CHAIN

If AND CHAIN is specified, a new (not aborted) transaction is immediately started with the same transaction characteristics (see [SET TRANSACTION](#)) as the just finished one. Otherwise, no new transaction is started.

Notes

Use [COMMIT](#) to successfully terminate a transaction.

Issuing ROLLBACK outside of a transaction block emits a warning and otherwise has no effect. ROLLBACK AND CHAIN outside of a transaction block is an error.

Examples

To abort all changes:

```
ROLLBACK;
```

Compatibility

The command ROLLBACK conforms to the SQL standard. The form ROLLBACK TRANSACTION is a Postgres Pro extension.

See Also

[BEGIN](#), [COMMIT](#), [ROLLBACK TO SAVEPOINT](#)

ROLLBACK PREPARED

ROLLBACK PREPARED — cancel a transaction that was earlier prepared for two-phase commit

Synopsis

```
ROLLBACK PREPARED transaction_id
```

Description

ROLLBACK PREPARED rolls back a transaction that is in prepared state.

Parameters

transaction_id

The transaction identifier of the transaction that is to be rolled back.

Notes

To roll back a prepared transaction, you must be either the same user that executed the transaction originally, or a superuser. But you do not have to be in the same session that executed the transaction.

This command cannot be executed inside a transaction block. The prepared transaction is rolled back immediately.

All currently available prepared transactions are listed in the [pg_prepared_xacts](#) system view.

Examples

Roll back the transaction identified by the transaction identifier `foobar`:

```
ROLLBACK PREPARED 'foobar';
```

Compatibility

ROLLBACK PREPARED is a Postgres Pro extension. It is intended for use by external transaction management systems, some of which are covered by standards (such as X/Open XA), but the SQL side of those systems is not standardized.

See Also

[PREPARE TRANSACTION](#), [COMMIT PREPARED](#)

ROLLBACK TO SAVEPOINT

ROLLBACK TO SAVEPOINT — roll back to a savepoint

Synopsis

```
ROLLBACK [ WORK | TRANSACTION ] TO [ SAVEPOINT ] savepoint_name
```

Description

Roll back all commands that were executed after the savepoint was established and then start a new subtransaction at the same transaction level. The savepoint remains valid and can be rolled back to again later, if needed.

ROLLBACK TO SAVEPOINT implicitly destroys all savepoints that were established after the named savepoint.

Parameters

savepoint_name

The savepoint to roll back to.

Notes

Use [RELEASE SAVEPOINT](#) to destroy a savepoint without discarding the effects of commands executed after it was established.

Specifying a savepoint name that has not been established is an error.

Cursors have somewhat non-transactional behavior with respect to savepoints. Any cursor that is opened inside a savepoint will be closed when the savepoint is rolled back. If a previously opened cursor is affected by a `FETCH` or `MOVE` command inside a savepoint that is later rolled back, the cursor remains at the position that `FETCH` left it pointing to (that is, the cursor motion caused by `FETCH` is not rolled back). Closing a cursor is not undone by rolling back, either. However, other side-effects caused by the cursor's query (such as side-effects of volatile functions called by the query) *are* rolled back if they occur during a savepoint that is later rolled back. A cursor whose execution causes a transaction to abort is put in a cannot-execute state, so while the transaction can be restored using `ROLLBACK TO SAVEPOINT`, the cursor can no longer be used.

Examples

To undo the effects of the commands executed after `my_savepoint` was established:

```
ROLLBACK TO SAVEPOINT my_savepoint;
```

Cursor positions are not affected by savepoint rollback:

```
BEGIN;

DECLARE foo CURSOR FOR SELECT 1 UNION SELECT 2;

SAVEPOINT foo;

FETCH 1 FROM foo;
  ?column?
-----
         1

ROLLBACK TO SAVEPOINT foo;
```



```
FETCH 1 FROM foo;  
?column?  
-----  
2
```

```
COMMIT;
```

Compatibility

The SQL standard specifies that the key word `SAVEPOINT` is mandatory, but Postgres Pro and Oracle allow it to be omitted. SQL allows only `WORK`, not `TRANSACTION`, as a noise word after `ROLLBACK`. Also, SQL has an optional clause `AND [NO] CHAIN` which is not currently supported by Postgres Pro. Otherwise, this command conforms to the SQL standard.

See Also

[BEGIN](#), [COMMIT](#), [RELEASE SAVEPOINT](#), [ROLLBACK](#), [SAVEPOINT](#)

SAVEPOINT

SAVEPOINT — define a new savepoint within the current transaction

Synopsis

```
SAVEPOINT savepoint_name
```

Description

SAVEPOINT establishes a new savepoint within the current transaction.

A savepoint is a special mark inside a transaction that allows all commands that are executed after it was established to be rolled back, restoring the transaction state to what it was at the time of the savepoint.

Parameters

savepoint_name

The name to give to the new savepoint. If savepoints with the same name already exist, they will be inaccessible until newer identically-named savepoints are released.

Notes

Use `ROLLBACK TO` to rollback to a savepoint. Use `RELEASE SAVEPOINT` to destroy a savepoint, keeping the effects of commands executed after it was established.

Savepoints can only be established when inside a transaction block. There can be multiple savepoints defined within a transaction.

Examples

To establish a savepoint and later undo the effects of all commands executed after it was established:

```
BEGIN;  
  INSERT INTO table1 VALUES (1);  
  SAVEPOINT my_savepoint;  
  INSERT INTO table1 VALUES (2);  
  ROLLBACK TO SAVEPOINT my_savepoint;  
  INSERT INTO table1 VALUES (3);  
COMMIT;
```

The above transaction will insert the values 1 and 3, but not 2.

To establish and later destroy a savepoint:

```
BEGIN;  
  INSERT INTO table1 VALUES (3);  
  SAVEPOINT my_savepoint;  
  INSERT INTO table1 VALUES (4);  
  RELEASE SAVEPOINT my_savepoint;  
COMMIT;
```

The above transaction will insert both 3 and 4.

To use a single savepoint name:

```
BEGIN;  
  INSERT INTO table1 VALUES (1);  
  SAVEPOINT my_savepoint;  
  INSERT INTO table1 VALUES (2);
```

```
SAVEPOINT my_savepoint;
INSERT INTO table1 VALUES (3);

-- rollback to the second savepoint
ROLLBACK TO SAVEPOINT my_savepoint;
SELECT * FROM table1;           -- shows rows 1 and 2

-- release the second savepoint
RELEASE SAVEPOINT my_savepoint;

-- rollback to the first savepoint
ROLLBACK TO SAVEPOINT my_savepoint;
SELECT * FROM table1;           -- shows only row 1
COMMIT;
```

The above transaction shows row 3 being rolled back first, then row 2.

Compatibility

SQL requires a savepoint to be destroyed automatically when another savepoint with the same name is established. In Postgres Pro, the old savepoint is kept, though only the more recent one will be used when rolling back or releasing. (Releasing the newer savepoint with `RELEASE SAVEPOINT` will cause the older one to again become accessible to `ROLLBACK TO SAVEPOINT` and `RELEASE SAVEPOINT`.) Otherwise, `SAVEPOINT` is fully SQL conforming.

See Also

[BEGIN](#), [COMMIT](#), [RELEASE SAVEPOINT](#), [ROLLBACK](#), [ROLLBACK TO SAVEPOINT](#)

SECURITY LABEL

SECURITY LABEL — define or change a security label applied to an object

Synopsis

```
SECURITY LABEL [ FOR provider ] ON
{
    TABLE object_name |
    COLUMN table_name.column_name |
    AGGREGATE aggregate_name ( aggregate_signature ) |
    DATABASE object_name |
    DOMAIN object_name |
    EVENT TRIGGER object_name |
    FOREIGN TABLE object_name |
    FUNCTION function_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    LARGE OBJECT large_object_oid |
    MATERIALIZED VIEW object_name |
    [ PROCEDURAL ] LANGUAGE object_name |
    PROCEDURE procedure_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    PUBLICATION object_name |
    PROFILE object_name |
    ROLE object_name |
    ROUTINE routine_name [ ( [ [ argmode ] [ argname ] argtype [, ...] ] ) ] |
    SCHEMA object_name |
    SEQUENCE object_name |
    SUBSCRIPTION object_name |
    TABLESPACE object_name |
    TYPE object_name |
    VIEW object_name
} IS { string_literal | NULL }

where aggregate_signature is:

* |
[ argmode ] [ argname ] argtype [ , ... ] |
[ [ argmode ] [ argname ] argtype [ , ... ] ] ORDER BY [ argmode ] [ argname ] argtype
[ , ... ]
```

Description

SECURITY LABEL applies a security label to a database object. An arbitrary number of security labels, one per label provider, can be associated with a given database object. Label providers are loadable modules which register themselves by using the function `register_label_provider`.

Note

`register_label_provider` is not an SQL function; it can only be called from C code loaded into the backend.

The label provider determines whether a given label is valid and whether it is permissible to assign that label to a given object. The meaning of a given label is likewise at the discretion of the label provider. Postgres Pro places no restrictions on whether or how a label provider must interpret security labels; it merely provides a mechanism for storing them. In practice, this facility is intended to allow integration with label-based mandatory access control (MAC) systems such as SELinux. Such systems make all

access control decisions based on object labels, rather than traditional discretionary access control (DAC) concepts such as users and groups.

Parameters

object_name
table_name.column_name
aggregate_name
function_name
procedure_name
routine_name

The name of the object to be labeled. Names of objects that reside in schemas (tables, functions, etc.) can be schema-qualified.

provider

The name of the provider with which this label is to be associated. The named provider must be loaded and must consent to the proposed labeling operation. If exactly one provider is loaded, the provider name may be omitted for brevity.

argmode

The mode of a function, procedure, or aggregate argument: IN, OUT, INOUT, or VARIADIC. If omitted, the default is IN. Note that SECURITY LABEL does not actually pay any attention to OUT arguments, since only the input arguments are needed to determine the function's identity. So it is sufficient to list the IN, INOUT, and VARIADIC arguments.

argname

The name of a function, procedure, or aggregate argument. Note that SECURITY LABEL does not actually pay any attention to argument names, since only the argument data types are needed to determine the function's identity.

argtype

The data type of a function, procedure, or aggregate argument.

large_object_oid

The OID of the large object.

PROCEDURAL

This is a noise word.

string_literal

The new setting of the security label, written as a string literal.

NULL

Write NULL to drop the security label.

Examples

The following example shows how the security label of a table could be set or changed:

```
SECURITY LABEL FOR selinux ON TABLE mytable IS 'system_u:object_r:sepgsql_table_t:s0';
```

To remove the label:

```
SECURITY LABEL FOR selinux ON TABLE mytable IS NULL;
```

Compatibility

There is no SECURITY LABEL command in the SQL standard.

SELECT

SELECT, TABLE, WITH — retrieve rows from a table or view

Synopsis

```
[ WITH [ RECURSIVE ] with_query [, ...] ]
SELECT [ ALL | DISTINCT [ ON ( expression [, ...] ) ] ]
    [ { * | expression [ [ AS ] output_name ] } [, ...] ]
    [ FROM from_item [, ...] ]
    [ WHERE condition ]
    [ GROUP BY [ ALL | DISTINCT ] grouping_element [, ...] ]
    [ HAVING condition ]
    [ WINDOW window_name AS ( window_definition ) [, ...] ]
    [ { UNION | INTERSECT | EXCEPT } [ ALL | DISTINCT ] select ]
    [ ORDER BY expression [ ASC | DESC | USING operator ] [ NULLS { FIRST | LAST } ]
[, ...] ]
    [ LIMIT { count | ALL } ]
    [ OFFSET start [ ROW | ROWS ] ]
    [ FETCH { FIRST | NEXT } [ count ] { ROW | ROWS } { ONLY | WITH TIES } ]
    [ FOR { UPDATE | NO KEY UPDATE | SHARE | KEY SHARE } [ OF from_reference [, ...] ]
[ NOWAIT | SKIP LOCKED ] [...]
```

where *from_item* can be one of:

```
    [ ONLY ] table_name [ * ] [ [ AS ] alias [ ( column_alias [, ...] ) ] ]
        [ TABLESAMPLE sampling_method ( argument [, ...] ) [ REPEATABLE ( seed
) ] ]
    [ LATERAL ] ( select ) [ [ AS ] alias [ ( column_alias [, ...] ) ] ]
    with_query_name [ [ AS ] alias [ ( column_alias [, ...] ) ] ]
    [ LATERAL ] function_name ( [ argument [, ...] ] )
        [ WITH ORDINALITY ] [ [ AS ] alias [ ( column_alias [, ...] ) ] ]
    [ LATERAL ] function_name ( [ argument [, ...] ] ) [ AS ] alias ( column_definition
[, ...] )
    [ LATERAL ] function_name ( [ argument [, ...] ] ) AS ( column_definition [, ...] )
    [ LATERAL ] ROWS FROM( function_name ( [ argument [, ...] ] ) [ AS
( column_definition [, ...] ) ] [, ...] )
        [ WITH ORDINALITY ] [ [ AS ] alias [ ( column_alias [, ...] ) ] ]
    from_item join_type from_item { ON join_condition | USING ( join_column [, ...] )
[ AS join_using_alias ] }
    from_item NATURAL join_type from_item
    from_item CROSS JOIN from_item
```

and *grouping_element* can be one of:

```
( )
expression
( expression [, ...] )
ROLLUP ( { expression | ( expression [, ...] ) } [, ...] )
CUBE ( { expression | ( expression [, ...] ) } [, ...] )
GROUPING SETS ( grouping_element [, ...] )
```

and *with_query* is:

```
    with_query_name [ ( column_name [, ...] ) ] AS [ [ NOT ] MATERIALIZED ] ( select
| values | insert | update | delete | merge )
```

```
[ SEARCH { BREADTH | DEPTH } FIRST BY column_name [, ...]
SET search_seq_col_name ]
[ CYCLE column_name [, ...] SET cycle_mark_col_name [ TO cycle_mark_value
DEFAULT cycle_mark_default ] USING cycle_path_col_name ]

TABLE [ ONLY ] table_name [ * ]
```

Description

SELECT retrieves rows from zero or more tables. The general processing of SELECT is as follows:

1. All queries in the WITH list are computed. These effectively serve as temporary tables that can be referenced in the FROM list. A WITH query that is referenced more than once in FROM is computed only once, unless specified otherwise with NOT MATERIALIZED. (See [WITH Clause](#) below.)
2. All elements in the FROM list are computed. (Each element in the FROM list is a real or virtual table.) If more than one element is specified in the FROM list, they are cross-joined together. (See [FROM Clause](#) below.)
3. If the WHERE clause is specified, all rows that do not satisfy the condition are eliminated from the output. (See [WHERE Clause](#) below.)
4. If the GROUP BY clause is specified, or if there are aggregate function calls, the output is combined into groups of rows that match on one or more values, and the results of aggregate functions are computed. If the HAVING clause is present, it eliminates groups that do not satisfy the given condition. (See [GROUP BY Clause](#) and [HAVING Clause](#) below.) Although query output columns are nominally computed in the next step, they can also be referenced (by name or ordinal number) in the GROUP BY clause.
5. The actual output rows are computed using the SELECT output expressions for each selected row or row group. (See [SELECT List](#) below.)
6. SELECT DISTINCT eliminates duplicate rows from the result. SELECT DISTINCT ON eliminates rows that match on all the specified expressions. SELECT ALL (the default) will return all candidate rows, including duplicates. (See [DISTINCT Clause](#) below.)
7. Using the operators UNION, INTERSECT, and EXCEPT, the output of more than one SELECT statement can be combined to form a single result set. The UNION operator returns all rows that are in one or both of the result sets. The INTERSECT operator returns all rows that are strictly in both result sets. The EXCEPT operator returns the rows that are in the first result set but not in the second. In all three cases, duplicate rows are eliminated unless ALL is specified. The noise word DISTINCT can be added to explicitly specify eliminating duplicate rows. Notice that DISTINCT is the default behavior here, even though ALL is the default for SELECT itself. (See [UNION Clause](#), [INTERSECT Clause](#), and [EXCEPT Clause](#) below.)
8. If the ORDER BY clause is specified, the returned rows are sorted in the specified order. If ORDER BY is not given, the rows are returned in whatever order the system finds fastest to produce. (See [ORDER BY Clause](#) below.)
9. If the LIMIT (or FETCH FIRST) or OFFSET clause is specified, the SELECT statement only returns a subset of the result rows. (See [LIMIT Clause](#) below.)
10. If FOR UPDATE, FOR NO KEY UPDATE, FOR SHARE or FOR KEY SHARE is specified, the SELECT statement locks the selected rows against concurrent updates. (See [The Locking Clause](#) below.)

You must have SELECT privilege on each column used in a SELECT command. The use of FOR NO KEY UPDATE, FOR UPDATE, FOR SHARE or FOR KEY SHARE requires UPDATE privilege as well (for at least one column of each table so selected).

Parameters

WITH Clause

The WITH clause allows you to specify one or more subqueries that can be referenced by name in the primary query. The subqueries effectively act as temporary tables or views for the duration of the primary

query. Each subquery can be a `SELECT`, `TABLE`, `VALUES`, `INSERT`, `UPDATE`, `DELETE`, or `MERGE` statement. When writing a data-modifying statement (`INSERT`, `UPDATE`, `DELETE`, or `MERGE`) in `WITH`, it is usual to include a `RETURNING` clause. It is the output of `RETURNING`, *not* the underlying table that the statement modifies, that forms the temporary table that is read by the primary query. If `RETURNING` is omitted, the statement is still executed, but it produces no output so it cannot be referenced as a table by the primary query.

A name (without schema qualification) must be specified for each `WITH` query. Optionally, a list of column names can be specified; if this is omitted, the column names are inferred from the subquery.

If `RECURSIVE` is specified, it allows a `SELECT` subquery to reference itself by name. Such a subquery must have the form

```
non_recursive_term UNION [ ALL | DISTINCT ] recursive_term
```

where the recursive self-reference must appear on the right-hand side of the `UNION`. Only one recursive self-reference is permitted per query. Recursive data-modifying statements are not supported, but you can use the results of a recursive `SELECT` query in a data-modifying statement. See [Section 7.8](#) for an example.

Another effect of `RECURSIVE` is that `WITH` queries need not be ordered: a query can reference another one that is later in the list. (However, circular references, or mutual recursion, are not implemented.) Without `RECURSIVE`, `WITH` queries can only reference sibling `WITH` queries that are earlier in the `WITH` list.

When there are multiple queries in the `WITH` clause, `RECURSIVE` should be written only once, immediately after `WITH`. It applies to all queries in the `WITH` clause, though it has no effect on queries that do not use recursion or forward references.

The optional `SEARCH` clause computes a *search sequence column* that can be used for ordering the results of a recursive query in either breadth-first or depth-first order. The supplied column name list specifies the row key that is to be used for keeping track of visited rows. A column named `search_seq_col_name` will be added to the result column list of the `WITH` query. This column can be ordered by in the outer query to achieve the respective ordering. See [Section 7.8.2.1](#) for examples.

The optional `CYCLE` clause is used to detect cycles in recursive queries. The supplied column name list specifies the row key that is to be used for keeping track of visited rows. A column named `cycle_mark_col_name` will be added to the result column list of the `WITH` query. This column will be set to `cycle_mark_value` when a cycle has been detected, else to `cycle_mark_default`. Furthermore, processing of the recursive union will stop when a cycle has been detected. `cycle_mark_value` and `cycle_mark_default` must be constants and they must be coercible to a common data type, and the data type must have an inequality operator. (The SQL standard requires that they be Boolean constants or character strings, but Postgres Pro does not require that.) By default, `TRUE` and `FALSE` (of type `boolean`) are used. Furthermore, a column named `cycle_path_col_name` will be added to the result column list of the `WITH` query. This column is used internally for tracking visited rows. See [Section 7.8.2.2](#) for examples.

Both the `SEARCH` and the `CYCLE` clause are only valid for recursive `WITH` queries. The `with_query` must be a `UNION` (or `UNION ALL`) of two `SELECT` (or equivalent) commands (no nested `UNIONS`). If both clauses are used, the column added by the `SEARCH` clause appears before the columns added by the `CYCLE` clause.

The primary query and the `WITH` queries are all (notionally) executed at the same time. This implies that the effects of a data-modifying statement in `WITH` cannot be seen from other parts of the query, other than by reading its `RETURNING` output. If two such data-modifying statements attempt to modify the same row, the results are unspecified.

A key property of `WITH` queries is that they are normally evaluated only once per execution of the primary query, even if the primary query refers to them more than once. In particular, data-modifying statements are guaranteed to be executed once and only once, regardless of whether the primary query reads all or any of their output.

However, a `WITH` query can be marked `NOT MATERIALIZED` to remove this guarantee. In that case, the `WITH` query can be folded into the primary query much as though it were a simple sub-`SELECT` in the

primary query's `FROM` clause. This results in duplicate computations if the primary query refers to that `WITH` query more than once; but if each such use requires only a few rows of the `WITH` query's total output, `NOT MATERIALIZED` can provide a net savings by allowing the queries to be optimized jointly. `NOT MATERIALIZED` is ignored if it is attached to a `WITH` query that is recursive or is not side-effect-free (i.e., is not a plain `SELECT` containing no volatile functions).

By default, a side-effect-free `WITH` query is folded into the primary query if it is used exactly once in the primary query's `FROM` clause. This allows joint optimization of the two query levels in situations where that should be semantically invisible. However, such folding can be prevented by marking the `WITH` query as `MATERIALIZED`. That might be useful, for example, if the `WITH` query is being used as an optimization fence to prevent the planner from choosing a bad plan. Postgres Pro versions before v12 never did such folding, so queries written for older versions might rely on `WITH` to act as an optimization fence.

See [Section 7.8](#) for additional information.

FROM Clause

The `FROM` clause specifies one or more source tables for the `SELECT`. If multiple sources are specified, the result is the Cartesian product (cross join) of all the sources. But usually qualification conditions are added (via `WHERE`) to restrict the returned rows to a small subset of the Cartesian product.

The `FROM` clause can contain the following elements:

table_name

The name (optionally schema-qualified) of an existing table or view. If `ONLY` is specified before the table name, only that table is scanned. If `ONLY` is not specified, the table and all its descendant tables (if any) are scanned. Optionally, `*` can be specified after the table name to explicitly indicate that descendant tables are included.

alias

A substitute name for the `FROM` item containing the alias. An alias is used for brevity or to eliminate ambiguity for self-joins (where the same table is scanned multiple times). When an alias is provided, it completely hides the actual name of the table or function; for example given `FROM foo AS f`, the remainder of the `SELECT` must refer to this `FROM` item as `f` not `foo`. If an alias is written, a column alias list can also be written to provide substitute names for one or more columns of the table.

`TABLESAMPLE sampling_method (argument [, ...]) [REPEATABLE (seed)]`

A `TABLESAMPLE` clause after a *table_name* indicates that the specified *sampling_method* should be used to retrieve a subset of the rows in that table. This sampling precedes the application of any other filters such as `WHERE` clauses. The standard Postgres Pro distribution includes two sampling methods, `BERNOULLI` and `SYSTEM`, and other sampling methods can be installed in the database via extensions.

The `BERNOULLI` and `SYSTEM` sampling methods each accept a single *argument* which is the fraction of the table to sample, expressed as a percentage between 0 and 100. This argument can be any real-valued expression. (Other sampling methods might accept more or different arguments.) These two methods each return a randomly-chosen sample of the table that will contain approximately the specified percentage of the table's rows. The `BERNOULLI` method scans the whole table and selects or ignores individual rows independently with the specified probability. The `SYSTEM` method does block-level sampling with each block having the specified chance of being selected; all rows in each selected block are returned. The `SYSTEM` method is significantly faster than the `BERNOULLI` method when small sampling percentages are specified, but it may return a less-random sample of the table as a result of clustering effects.

The optional `REPEATABLE` clause specifies a *seed* number or expression to use for generating random numbers within the sampling method. The seed value can be any non-null floating-point value. Two queries that specify the same seed and *argument* values will select the same sample of the table, if the table has not been changed meanwhile. But different seed values will usually produce different samples. If `REPEATABLE` is not given then a new random sample is selected for each query, based

upon a system-generated seed. Note that some add-on sampling methods do not accept `REPEATABLE`, and will always produce new samples on each use.

select

A sub-`SELECT` can appear in the `FROM` clause. This acts as though its output were created as a temporary table for the duration of this single `SELECT` command. Note that the sub-`SELECT` must be surrounded by parentheses, and an alias can be provided in the same way as for a table. A `VALUES` command can also be used here.

with_query_name

A `WITH` query is referenced by writing its name, just as though the query's name were a table name. (In fact, the `WITH` query hides any real table of the same name for the purposes of the primary query. If necessary, you can refer to a real table of the same name by schema-qualifying the table's name.) An alias can be provided in the same way as for a table.

function_name

Function calls can appear in the `FROM` clause. (This is especially useful for functions that return result sets, but any function can be used.) This acts as though the function's output were created as a temporary table for the duration of this single `SELECT` command. If the function's result type is composite (including the case of a function with multiple `OUT` parameters), each attribute becomes a separate column in the implicit table.

When the optional `WITH ORDINALITY` clause is added to the function call, an additional column of type `bigint` will be appended to the function's result column(s). This column numbers the rows of the function's result set, starting from 1. By default, this column is named `ordinality`.

An alias can be provided in the same way as for a table. If an alias is written, a column alias list can also be written to provide substitute names for one or more attributes of the function's composite return type, including the ordinality column if present.

Multiple function calls can be combined into a single `FROM`-clause item by surrounding them with `ROWS FROM(...)`. The output of such an item is the concatenation of the first row from each function, then the second row from each function, etc. If some of the functions produce fewer rows than others, null values are substituted for the missing data, so that the total number of rows returned is always the same as for the function that produced the most rows.

If the function has been defined as returning the `record` data type, then an alias or the key word `AS` must be present, followed by a column definition list in the form (`column_name data_type` [, ...]). The column definition list must match the actual number and types of columns returned by the function.

When using the `ROWS FROM(...)` syntax, if one of the functions requires a column definition list, it's preferred to put the column definition list after the function call inside `ROWS FROM(...)`. A column definition list can be placed after the `ROWS FROM(...)` construct only if there's just a single function and no `WITH ORDINALITY` clause.

To use `ORDINALITY` together with a column definition list, you must use the `ROWS FROM(...)` syntax and put the column definition list inside `ROWS FROM(...)`.

join_type

One of

- [`INNER`] `JOIN`
- `LEFT` [`OUTER`] `JOIN`
- `RIGHT` [`OUTER`] `JOIN`
- `FULL` [`OUTER`] `JOIN`

For the `INNER` and `OUTER` join types, a join condition must be specified, namely exactly one of `ON join_condition`, `USING (join_column [, ...])`, or `NATURAL`. See below for the meaning.

A `JOIN` clause combines two `FROM` items, which for convenience we will refer to as “tables”, though in reality they can be any type of `FROM` item. Use parentheses if necessary to determine the order of nesting. In the absence of parentheses, `JOINS` nest left-to-right. In any case `JOIN` binds more tightly than the commas separating `FROM`-list items. All the `JOIN` options are just a notational convenience, since they do nothing you couldn't do with plain `FROM` and `WHERE`.

`LEFT OUTER JOIN` returns all rows in the qualified Cartesian product (i.e., all combined rows that pass its join condition), plus one copy of each row in the left-hand table for which there was no right-hand row that passed the join condition. This left-hand row is extended to the full width of the joined table by inserting null values for the right-hand columns. Note that only the `JOIN` clause's own condition is considered while deciding which rows have matches. Outer conditions are applied afterwards.

Conversely, `RIGHT OUTER JOIN` returns all the joined rows, plus one row for each unmatched right-hand row (extended with nulls on the left). This is just a notational convenience, since you could convert it to a `LEFT OUTER JOIN` by switching the left and right tables.

`FULL OUTER JOIN` returns all the joined rows, plus one row for each unmatched left-hand row (extended with nulls on the right), plus one row for each unmatched right-hand row (extended with nulls on the left).

`ON join_condition`

join_condition is an expression resulting in a value of type `boolean` (similar to a `WHERE` clause) that specifies which rows in a join are considered to match.

`USING (join_column [, ...]) [AS join_using_alias]`

A clause of the form `USING (a, b, ...)` is shorthand for `ON left_table.a = right_table.a AND left_table.b = right_table.b` Also, `USING` implies that only one of each pair of equivalent columns will be included in the join output, not both.

If a *join_using_alias* name is specified, it provides a table alias for the join columns. Only the join columns listed in the `USING` clause are addressable by this name. Unlike a regular *alias*, this does not hide the names of the joined tables from the rest of the query. Also unlike a regular *alias*, you cannot write a column alias list — the output names of the join columns are the same as they appear in the `USING` list.

`NATURAL`

`NATURAL` is shorthand for a `USING` list that mentions all columns in the two tables that have matching names. If there are no common column names, `NATURAL` is equivalent to `ON TRUE`.

`CROSS JOIN`

`CROSS JOIN` is equivalent to `INNER JOIN ON (TRUE)`, that is, no rows are removed by qualification. They produce a simple Cartesian product, the same result as you get from listing the two tables at the top level of `FROM`, but restricted by the join condition (if any).

`LATERAL`

The `LATERAL` key word can precede a sub-`SELECT FROM` item. This allows the sub-`SELECT` to refer to columns of `FROM` items that appear before it in the `FROM` list. (Without `LATERAL`, each sub-`SELECT` is evaluated independently and so cannot cross-reference any other `FROM` item.)

`LATERAL` can also precede a function-call `FROM` item, but in this case it is a noise word, because the function expression can refer to earlier `FROM` items in any case.

A `LATERAL` item can appear at top level in the `FROM` list, or within a `JOIN` tree. In the latter case it can also refer to any items that are on the left-hand side of a `JOIN` that it is on the right-hand side of.

When a `FROM` item contains `LATERAL` cross-references, evaluation proceeds as follows: for each row of the `FROM` item providing the cross-referenced column(s), or set of rows of multiple `FROM` items providing the columns, the `LATERAL` item is evaluated using that row or row set's values of the columns. The resulting row(s) are joined as usual with the rows they were computed from. This is repeated for each row or set of rows from the column source table(s).

The column source table(s) must be `INNER` or `LEFT` joined to the `LATERAL` item, else there would not be a well-defined set of rows from which to compute each set of rows for the `LATERAL` item. Thus, although a construct such as `X RIGHT JOIN LATERAL Y` is syntactically valid, it is not actually allowed for `Y` to reference `X`.

WHERE Clause

The optional `WHERE` clause has the general form

```
WHERE condition
```

where *condition* is any expression that evaluates to a result of type `boolean`. Any row that does not satisfy this condition will be eliminated from the output. A row satisfies the condition if it returns true when the actual row values are substituted for any variable references.

GROUP BY Clause

The optional `GROUP BY` clause has the general form

```
GROUP BY [ ALL | DISTINCT ] grouping_element [, ...]
```

`GROUP BY` will condense into a single row all selected rows that share the same values for the grouped expressions. An *expression* used inside a *grouping_element* can be an input column name, or the name or ordinal number of an output column (`SELECT` list item), or an arbitrary expression formed from input-column values. In case of ambiguity, a `GROUP BY` name will be interpreted as an input-column name rather than an output column name.

If any of `GROUPING SETS`, `ROLLUP` or `CUBE` are present as grouping elements, then the `GROUP BY` clause as a whole defines some number of independent *grouping sets*. The effect of this is equivalent to constructing a `UNION ALL` between subqueries with the individual grouping sets as their `GROUP BY` clauses. The optional `DISTINCT` clause removes duplicate sets before processing; it does *not* transform the `UNION ALL` into a `UNION DISTINCT`. For further details on the handling of grouping sets see [Section 7.2.4](#).

Aggregate functions, if any are used, are computed across all rows making up each group, producing a separate value for each group. (If there are aggregate functions but no `GROUP BY` clause, the query is treated as having a single group comprising all the selected rows.) The set of rows fed to each aggregate function can be further filtered by attaching a `FILTER` clause to the aggregate function call; see [Section 4.2.7](#) for more information. When a `FILTER` clause is present, only those rows matching it are included in the input to that aggregate function.

When `GROUP BY` is present, or any aggregate functions are present, it is not valid for the `SELECT` list expressions to refer to ungrouped columns except within aggregate functions or when the ungrouped column is functionally dependent on the grouped columns, since there would otherwise be more than one possible value to return for an ungrouped column. A functional dependency exists if the grouped columns (or a subset thereof) are the primary key of the table containing the ungrouped column.

Keep in mind that all aggregate functions are evaluated before evaluating any “scalar” expressions in the `HAVING` clause or `SELECT` list. This means that, for example, a `CASE` expression cannot be used to skip evaluation of an aggregate function; see [Section 4.2.14](#).

Currently, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE` and `FOR KEY SHARE` cannot be specified with `GROUP BY`.

HAVING Clause

The optional `HAVING` clause has the general form

HAVING *condition*

where *condition* is the same as specified for the WHERE clause.

HAVING eliminates group rows that do not satisfy the condition. HAVING is different from WHERE: WHERE filters individual rows before the application of GROUP BY, while HAVING filters group rows created by GROUP BY. Each column referenced in *condition* must unambiguously reference a grouping column, unless the reference appears within an aggregate function or the ungrouped column is functionally dependent on the grouping columns.

The presence of HAVING turns a query into a grouped query even if there is no GROUP BY clause. This is the same as what happens when the query contains aggregate functions but no GROUP BY clause. All the selected rows are considered to form a single group, and the SELECT list and HAVING clause can only reference table columns from within aggregate functions. Such a query will emit a single row if the HAVING condition is true, zero rows if it is not true.

Currently, FOR NO KEY UPDATE, FOR UPDATE, FOR SHARE and FOR KEY SHARE cannot be specified with HAVING.

WINDOW Clause

The optional WINDOW clause has the general form

```
WINDOW window_name AS ( window_definition ) [, ...]
```

where *window_name* is a name that can be referenced from OVER clauses or subsequent window definitions, and *window_definition* is

```
[ existing_window_name ]  
[ PARTITION BY expression [, ...] ]  
[ ORDER BY expression [ ASC | DESC | USING operator ] [ NULLS { FIRST | LAST } ]  
  [, ...] ]  
[ frame_clause ]
```

If an *existing_window_name* is specified it must refer to an earlier entry in the WINDOW list; the new window copies its partitioning clause from that entry, as well as its ordering clause if any. In this case the new window cannot specify its own PARTITION BY clause, and it can specify ORDER BY only if the copied window does not have one. The new window always uses its own frame clause; the copied window must not specify a frame clause.

The elements of the PARTITION BY list are interpreted in much the same fashion as elements of a GROUP BY clause, except that they are always simple expressions and never the name or number of an output column. Another difference is that these expressions can contain aggregate function calls, which are not allowed in a regular GROUP BY clause. They are allowed here because windowing occurs after grouping and aggregation.

Similarly, the elements of the ORDER BY list are interpreted in much the same fashion as elements of a statement-level ORDER BY clause, except that the expressions are always taken as simple expressions and never the name or number of an output column.

The optional *frame_clause* defines the *window frame* for window functions that depend on the frame (not all do). The window frame is a set of related rows for each row of the query (called the *current row*). The *frame_clause* can be one of

```
{ RANGE | ROWS | GROUPS } frame_start [ frame_exclusion ]  
{ RANGE | ROWS | GROUPS } BETWEEN frame_start AND frame_end [ frame_exclusion ]
```

where *frame_start* and *frame_end* can be one of

```
UNBOUNDED PRECEDING  
offset PRECEDING  
CURRENT ROW  
offset FOLLOWING
```

UNBOUNDED FOLLOWING

and *frame_exclusion* can be one of

EXCLUDE CURRENT ROW
EXCLUDE GROUP
EXCLUDE TIES
EXCLUDE NO OTHERS

If *frame_end* is omitted it defaults to CURRENT ROW. Restrictions are that *frame_start* cannot be UNBOUNDED FOLLOWING, *frame_end* cannot be UNBOUNDED PRECEDING, and the *frame_end* choice cannot appear earlier in the above list of *frame_start* and *frame_end* options than the *frame_start* choice does — for example RANGE BETWEEN CURRENT ROW AND *offset* PRECEDING is not allowed.

The default framing option is RANGE UNBOUNDED PRECEDING, which is the same as RANGE BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW; it sets the frame to be all rows from the partition start up through the current row's last *peer* (a row that the window's ORDER BY clause considers equivalent to the current row; all rows are peers if there is no ORDER BY). In general, UNBOUNDED PRECEDING means that the frame starts with the first row of the partition, and similarly UNBOUNDED FOLLOWING means that the frame ends with the last row of the partition, regardless of RANGE, ROWS or GROUPS mode. In ROWS mode, CURRENT ROW means that the frame starts or ends with the current row; but in RANGE or GROUPS mode it means that the frame starts or ends with the current row's first or last peer in the ORDER BY ordering. The *offset* PRECEDING and *offset* FOLLOWING options vary in meaning depending on the frame mode. In ROWS mode, the *offset* is an integer indicating that the frame starts or ends that many rows before or after the current row. In GROUPS mode, the *offset* is an integer indicating that the frame starts or ends that many peer groups before or after the current row's peer group, where a *peer group* is a group of rows that are equivalent according to the window's ORDER BY clause. In RANGE mode, use of an *offset* option requires that there be exactly one ORDER BY column in the window definition. Then the frame contains those rows whose ordering column value is no more than *offset* less than (for PRECEDING) or more than (for FOLLOWING) the current row's ordering column value. In these cases the data type of the *offset* expression depends on the data type of the ordering column. For numeric ordering columns it is typically of the same type as the ordering column, but for datetime ordering columns it is an interval. In all these cases, the value of the *offset* must be non-null and non-negative. Also, while the *offset* does not have to be a simple constant, it cannot contain variables, aggregate functions, or window functions.

The *frame_exclusion* option allows rows around the current row to be excluded from the frame, even if they would be included according to the frame start and frame end options. EXCLUDE CURRENT ROW excludes the current row from the frame. EXCLUDE GROUP excludes the current row and its ordering peers from the frame. EXCLUDE TIES excludes any peers of the current row from the frame, but not the current row itself. EXCLUDE NO OTHERS simply specifies explicitly the default behavior of not excluding the current row or its peers.

Beware that the ROWS mode can produce unpredictable results if the ORDER BY ordering does not order the rows uniquely. The RANGE and GROUPS modes are designed to ensure that rows that are peers in the ORDER BY ordering are treated alike: all rows of a given peer group will be in the frame or excluded from it.

The purpose of a WINDOW clause is to specify the behavior of *window functions* appearing in the query's SELECT list or ORDER BY clause. These functions can reference the WINDOW clause entries by name in their OVER clauses. A WINDOW clause entry does not have to be referenced anywhere, however; if it is not used in the query it is simply ignored. It is possible to use window functions without any WINDOW clause at all, since a window function call can specify its window definition directly in its OVER clause. However, the WINDOW clause saves typing when the same window definition is needed for more than one window function.

Currently, FOR NO KEY UPDATE, FOR UPDATE, FOR SHARE and FOR KEY SHARE cannot be specified with WINDOW.

Window functions are described in detail in [Section 3.5](#), [Section 4.2.8](#), and [Section 7.2.5](#).

SELECT List

The `SELECT` list (between the key words `SELECT` and `FROM`) specifies expressions that form the output rows of the `SELECT` statement. The expressions can (and usually do) refer to columns computed in the `FROM` clause.

Just as in a table, every output column of a `SELECT` has a name. In a simple `SELECT` this name is just used to label the column for display, but when the `SELECT` is a sub-query of a larger query, the name is seen by the larger query as the column name of the virtual table produced by the sub-query. To specify the name to use for an output column, write `AS output_name` after the column's expression. (You can omit `AS`, but only if the desired output name does not match any Postgres Pro keyword (see [Appendix C](#)). For protection against possible future keyword additions, it is recommended that you always either write `AS` or double-quote the output name.) If you do not specify a column name, a name is chosen automatically by Postgres Pro. If the column's expression is a simple column reference then the chosen name is the same as that column's name. In more complex cases a function or type name may be used, or the system may fall back on a generated name such as `?column?`.

An output column's name can be used to refer to the column's value in `ORDER BY` and `GROUP BY` clauses, but not in the `WHERE` or `HAVING` clauses; there you must write out the expression instead.

Instead of an expression, `*` can be written in the output list as a shorthand for all the columns of the selected rows. Also, you can write `table_name.*` as a shorthand for the columns coming from just that table. In these cases it is not possible to specify new names with `AS`; the output column names will be the same as the table columns' names.

According to the SQL standard, the expressions in the output list should be computed before applying `DISTINCT`, `ORDER BY`, or `LIMIT`. This is obviously necessary when using `DISTINCT`, since otherwise it's not clear what values are being made distinct. However, in many cases it is convenient if output expressions are computed after `ORDER BY` and `LIMIT`; particularly if the output list contains any volatile or expensive functions. With that behavior, the order of function evaluations is more intuitive and there will not be evaluations corresponding to rows that never appear in the output. Postgres Pro will effectively evaluate output expressions after sorting and limiting, so long as those expressions are not referenced in `DISTINCT`, `ORDER BY` or `GROUP BY`. (As a counterexample, `SELECT f(x) FROM tab ORDER BY 1` clearly must evaluate `f(x)` before sorting.) Output expressions that contain set-returning functions are effectively evaluated after sorting and before limiting, so that `LIMIT` will act to cut off the output from a set-returning function.

Note

Postgres Pro versions before 9.6 did not provide any guarantees about the timing of evaluation of output expressions versus sorting and limiting; it depended on the form of the chosen query plan.

DISTINCT Clause

If `SELECT DISTINCT` is specified, all duplicate rows are removed from the result set (one row is kept from each group of duplicates). `SELECT ALL` specifies the opposite: all rows are kept; that is the default.

`SELECT DISTINCT ON (expression [, ...])` keeps only the first row of each set of rows where the given expressions evaluate to equal. The `DISTINCT ON` expressions are interpreted using the same rules as for `ORDER BY` (see above). Note that the “first row” of each set is unpredictable unless `ORDER BY` is used to ensure that the desired row appears first. For example:

```
SELECT DISTINCT ON (location) location, time, report
FROM weather_reports
ORDER BY location, time DESC;
```

retrieves the most recent weather report for each location. But if we had not used `ORDER BY` to force descending order of time values for each location, we'd have gotten a report from an unpredictable time for each location.

The `DISTINCT ON` expression(s) must match the leftmost `ORDER BY` expression(s). The `ORDER BY` clause will normally contain additional expression(s) that determine the desired precedence of rows within each `DISTINCT ON` group.

Currently, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE` and `FOR KEY SHARE` cannot be specified with `DISTINCT`.

UNION Clause

The `UNION` clause has this general form:

```
select_statement UNION [ ALL | DISTINCT ] select_statement
```

select_statement is any `SELECT` statement without an `ORDER BY`, `LIMIT`, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE`, or `FOR KEY SHARE` clause. (`ORDER BY` and `LIMIT` can be attached to a subexpression if it is enclosed in parentheses. Without parentheses, these clauses will be taken to apply to the result of the `UNION`, not to its right-hand input expression.)

The `UNION` operator computes the set union of the rows returned by the involved `SELECT` statements. A row is in the set union of two result sets if it appears in at least one of the result sets. The two `SELECT` statements that represent the direct operands of the `UNION` must produce the same number of columns, and corresponding columns must be of compatible data types.

The result of `UNION` does not contain any duplicate rows unless the `ALL` option is specified. `ALL` prevents elimination of duplicates. (Therefore, `UNION ALL` is usually significantly quicker than `UNION`; use `ALL` when you can.) `DISTINCT` can be written to explicitly specify the default behavior of eliminating duplicate rows.

Multiple `UNION` operators in the same `SELECT` statement are evaluated left to right, unless otherwise indicated by parentheses.

Currently, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE` and `FOR KEY SHARE` cannot be specified either for a `UNION` result or for any input of a `UNION`.

INTERSECT Clause

The `INTERSECT` clause has this general form:

```
select_statement INTERSECT [ ALL | DISTINCT ] select_statement
```

select_statement is any `SELECT` statement without an `ORDER BY`, `LIMIT`, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE`, or `FOR KEY SHARE` clause.

The `INTERSECT` operator computes the set intersection of the rows returned by the involved `SELECT` statements. A row is in the intersection of two result sets if it appears in both result sets.

The result of `INTERSECT` does not contain any duplicate rows unless the `ALL` option is specified. With `ALL`, a row that has m duplicates in the left table and n duplicates in the right table will appear $\min(m, n)$ times in the result set. `DISTINCT` can be written to explicitly specify the default behavior of eliminating duplicate rows.

Multiple `INTERSECT` operators in the same `SELECT` statement are evaluated left to right, unless parentheses dictate otherwise. `INTERSECT` binds more tightly than `UNION`. That is, `A UNION B INTERSECT C` will be read as `A UNION (B INTERSECT C)`.

Currently, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE` and `FOR KEY SHARE` cannot be specified either for an `INTERSECT` result or for any input of an `INTERSECT`.

EXCEPT Clause

The `EXCEPT` clause has this general form:

```
select_statement EXCEPT [ ALL | DISTINCT ] select_statement
```

select_statement is any `SELECT` statement without an `ORDER BY`, `LIMIT`, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE`, or `FOR KEY SHARE` clause.

The `EXCEPT` operator computes the set of rows that are in the result of the left `SELECT` statement but not in the result of the right one.

The result of `EXCEPT` does not contain any duplicate rows unless the `ALL` option is specified. With `ALL`, a row that has m duplicates in the left table and n duplicates in the right table will appear $\max(m-n, 0)$ times in the result set. `DISTINCT` can be written to explicitly specify the default behavior of eliminating duplicate rows.

Multiple `EXCEPT` operators in the same `SELECT` statement are evaluated left to right, unless parentheses dictate otherwise. `EXCEPT` binds at the same level as `UNION`.

Currently, `FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE` and `FOR KEY SHARE` cannot be specified either for an `EXCEPT` result or for any input of an `EXCEPT`.

ORDER BY Clause

The optional `ORDER BY` clause has this general form:

```
ORDER BY expression [ ASC | DESC | USING operator ] [ NULLS { FIRST | LAST } ] [, ...]
```

The `ORDER BY` clause causes the result rows to be sorted according to the specified expression(s). If two rows are equal according to the leftmost expression, they are compared according to the next expression and so on. If they are equal according to all specified expressions, they are returned in an implementation-dependent order.

Each *expression* can be the name or ordinal number of an output column (`SELECT` list item), or it can be an arbitrary expression formed from input-column values.

The ordinal number refers to the ordinal (left-to-right) position of the output column. This feature makes it possible to define an ordering on the basis of a column that does not have a unique name. This is never absolutely necessary because it is always possible to assign a name to an output column using the `AS` clause.

It is also possible to use arbitrary expressions in the `ORDER BY` clause, including columns that do not appear in the `SELECT` output list. Thus the following statement is valid:

```
SELECT name FROM distributors ORDER BY code;
```

A limitation of this feature is that an `ORDER BY` clause applying to the result of a `UNION`, `INTERSECT`, or `EXCEPT` clause can only specify an output column name or number, not an expression.

If an `ORDER BY` expression is a simple name that matches both an output column name and an input column name, `ORDER BY` will interpret it as the output column name. This is the opposite of the choice that `GROUP BY` will make in the same situation. This inconsistency is made to be compatible with the SQL standard.

Optionally one can add the key word `ASC` (ascending) or `DESC` (descending) after any expression in the `ORDER BY` clause. If not specified, `ASC` is assumed by default. Alternatively, a specific ordering operator name can be specified in the `USING` clause. An ordering operator must be a less-than or greater-than member of some B-tree operator family. `ASC` is usually equivalent to `USING <` and `DESC` is usually equivalent to `USING >`. (But the creator of a user-defined data type can define exactly what the default sort ordering is, and it might correspond to operators with other names.)

If `NULLS LAST` is specified, null values sort after all non-null values; if `NULLS FIRST` is specified, null values sort before all non-null values. If neither is specified, the default behavior is `NULLS LAST` when `ASC` is specified or implied, and `NULLS FIRST` when `DESC` is specified (thus, the default is to act as though nulls are larger than non-nulls). When `USING` is specified, the default nulls ordering depends on whether the operator is a less-than or greater-than operator.

Note that ordering options apply only to the expression they follow; for example `ORDER BY x, y DESC` does not mean the same thing as `ORDER BY x DESC, y DESC`.

Character-string data is sorted according to the collation that applies to the column being sorted. That can be overridden at need by including a `COLLATE` clause in the *expression*, for example `ORDER BY mycolumn COLLATE "en_US"`. For more information see [Section 4.2.10](#) and [Section 23.2](#).

LIMIT Clause

The `LIMIT` clause consists of two independent sub-clauses:

```
LIMIT { count | ALL }  
OFFSET start
```

The parameter *count* specifies the maximum number of rows to return, while *start* specifies the number of rows to skip before starting to return rows. When both are specified, *start* rows are skipped before starting to count the *count* rows to be returned.

If the *count* expression evaluates to `NULL`, it is treated as `LIMIT ALL`, i.e., no limit. If *start* evaluates to `NULL`, it is treated the same as `OFFSET 0`.

SQL:2008 introduced a different syntax to achieve the same result, which Postgres Pro also supports. It is:

```
OFFSET start { ROW | ROWS }  
FETCH { FIRST | NEXT } [ count ] { ROW | ROWS } { ONLY | WITH TIES }
```

In this syntax, the *start* or *count* value is required by the standard to be a literal constant, a parameter, or a variable name; as a Postgres Pro extension, other expressions are allowed, but will generally need to be enclosed in parentheses to avoid ambiguity. If *count* is omitted in a `FETCH` clause, it defaults to 1. The `WITH TIES` option is used to return any additional rows that tie for the last place in the result set according to the `ORDER BY` clause; `ORDER BY` is mandatory in this case, and `SKIP LOCKED` is not allowed. `ROW` and `ROWS` as well as `FIRST` and `NEXT` are noise words that don't influence the effects of these clauses. According to the standard, the `OFFSET` clause must come before the `FETCH` clause if both are present; but Postgres Pro is laxer and allows either order.

When using `LIMIT`, it is a good idea to use an `ORDER BY` clause that constrains the result rows into a unique order. Otherwise you will get an unpredictable subset of the query's rows — you might be asking for the tenth through twentieth rows, but tenth through twentieth in what ordering? You don't know what ordering unless you specify `ORDER BY`.

The query planner takes `LIMIT` into account when generating a query plan, so you are very likely to get different plans (yielding different row orders) depending on what you use for `LIMIT` and `OFFSET`. Thus, using different `LIMIT/OFFSET` values to select different subsets of a query result *will give inconsistent results* unless you enforce a predictable result ordering with `ORDER BY`. This is not a bug; it is an inherent consequence of the fact that SQL does not promise to deliver the results of a query in any particular order unless `ORDER BY` is used to constrain the order.

It is even possible for repeated executions of the same `LIMIT` query to return different subsets of the rows of a table, if there is not an `ORDER BY` to enforce selection of a deterministic subset. Again, this is not a bug; determinism of the results is simply not guaranteed in such a case.

The Locking Clause

`FOR UPDATE`, `FOR NO KEY UPDATE`, `FOR SHARE` and `FOR KEY SHARE` are *locking clauses*; they affect how `SELECT` locks rows as they are obtained from the table.

The locking clause has the general form

```
FOR lock_strength [ OF from_reference [, ...] ] [ NOWAIT | SKIP LOCKED ]
```

where *lock_strength* can be one of

```
UPDATE  
NO KEY UPDATE  
SHARE  
KEY SHARE
```

from_reference must be a table *alias* or non-hidden *table_name* referenced in the `FROM` clause. For more information on each row-level lock mode, refer to [Section 13.3.2](#).

To prevent the operation from waiting for other transactions to commit, use either the `NOWAIT` or `SKIP LOCKED` option. With `NOWAIT`, the statement reports an error, rather than waiting, if a selected row cannot be locked immediately. With `SKIP LOCKED`, any selected rows that cannot be immediately locked are skipped. Skipping locked rows provides an inconsistent view of the data, so this is not suitable for general purpose work, but can be used to avoid lock contention with multiple consumers accessing a queue-like table. Note that `NOWAIT` and `SKIP LOCKED` apply only to the row-level lock(s) — the required `ROW SHARE` table-level lock is still taken in the ordinary way (see [Chapter 13](#)). You can use `LOCK` with the `NOWAIT` option first, if you need to acquire the table-level lock without waiting.

If specific tables are named in a locking clause, then only rows coming from those tables are locked; any other tables used in the `SELECT` are simply read as usual. A locking clause without a table list affects all tables used in the statement. If a locking clause is applied to a view or sub-query, it affects all tables used in the view or sub-query. However, these clauses do not apply to `WITH` queries referenced by the primary query. If you want row locking to occur within a `WITH` query, specify a locking clause within the `WITH` query.

Multiple locking clauses can be written if it is necessary to specify different locking behavior for different tables. If the same table is mentioned (or implicitly affected) by more than one locking clause, then it is processed as if it was only specified by the strongest one. Similarly, a table is processed as `NOWAIT` if that is specified in any of the clauses affecting it. Otherwise, it is processed as `SKIP LOCKED` if that is specified in any of the clauses affecting it.

The locking clauses cannot be used in contexts where returned rows cannot be clearly identified with individual table rows; for example they cannot be used with aggregation.

When a locking clause appears at the top level of a `SELECT` query, the rows that are locked are exactly those that are returned by the query; in the case of a join query, the rows locked are those that contribute to returned join rows. In addition, rows that satisfied the query conditions as of the query snapshot will be locked, although they will not be returned if they were updated after the snapshot and no longer satisfy the query conditions. If a `LIMIT` is used, locking stops once enough rows have been returned to satisfy the limit (but note that rows skipped over by `OFFSET` will get locked). Similarly, if a locking clause is used in a cursor's query, only rows actually fetched or stepped past by the cursor will be locked.

When a locking clause appears in a sub-`SELECT`, the rows locked are those returned to the outer query by the sub-query. This might involve fewer rows than inspection of the sub-query alone would suggest, since conditions from the outer query might be used to optimize execution of the sub-query. For example,

```
SELECT * FROM (SELECT * FROM mytable FOR UPDATE) ss WHERE col1 = 5;
```

will lock only rows having `col1 = 5`, even though that condition is not textually within the sub-query.

Previous releases failed to preserve a lock which is upgraded by a later savepoint. For example, this code:

```
BEGIN;
SELECT * FROM mytable WHERE key = 1 FOR UPDATE;
SAVEPOINT s;
UPDATE mytable SET ... WHERE key = 1;
ROLLBACK TO s;
```

would fail to preserve the `FOR UPDATE` lock after the `ROLLBACK TO`. This has been fixed in release 9.3.

Caution

It is possible for a `SELECT` command running at the `READ COMMITTED` transaction isolation level and using `ORDER BY` and a locking clause to return rows out of order. This is because `ORDER BY` is applied first. The command sorts the result, but might then block trying to obtain a lock on one or more of the rows. Once the `SELECT` unblocks, some of the ordering column values might have been

modified, leading to those rows appearing to be out of order (though they are in order in terms of the original column values). This can be worked around at need by placing the `FOR UPDATE/SHARE` clause in a sub-query, for example

```
SELECT * FROM (SELECT * FROM mytable FOR UPDATE) ss ORDER BY column1;
```

Note that this will result in locking all rows of `mytable`, whereas `FOR UPDATE` at the top level would lock only the actually returned rows. This can make for a significant performance difference, particularly if the `ORDER BY` is combined with `LIMIT` or other restrictions. So this technique is recommended only if concurrent updates of the ordering columns are expected and a strictly sorted result is required.

At the `REPEATABLE READ` or `SERIALIZABLE` transaction isolation level this would cause a serialization failure (with an `SQLSTATE` of `'40001'`), so there is no possibility of receiving rows out of order under these isolation levels.

TABLE Command

The command

```
TABLE name
```

is equivalent to

```
SELECT * FROM name
```

It can be used as a top-level command or as a space-saving syntax variant in parts of complex queries. Only the `WITH`, `UNION`, `INTERSECT`, `EXCEPT`, `ORDER BY`, `LIMIT`, `OFFSET`, `FETCH` and `FOR` locking clauses can be used with `TABLE`; the `WHERE` clause and any form of aggregation cannot be used.

Examples

To join the table `films` with the table `distributors`:

```
SELECT f.title, f.did, d.name, f.date_prod, f.kind
FROM distributors d JOIN films f USING (did);
```

title	did	name	date_prod	kind
The Third Man	101	British Lion	1949-12-23	Drama
The African Queen	101	British Lion	1951-08-11	Romantic
...				

To sum the column `len` of all films and group the results by `kind`:

```
SELECT kind, sum(len) AS total FROM films GROUP BY kind;
```

kind	total
Action	07:34
Comedy	02:58
Drama	14:28
Musical	06:42
Romantic	04:38

To sum the column `len` of all films, group the results by `kind` and show those group totals that are less than 5 hours:

```
SELECT kind, sum(len) AS total
FROM films
GROUP BY kind
HAVING sum(len) < interval '5 hours';
```

SELECT

kind	total
Comedy	02:58
Romantic	04:38

The following two examples are identical ways of sorting the individual results according to the contents of the second column (name):

```
SELECT * FROM distributors ORDER BY name;
SELECT * FROM distributors ORDER BY 2;
```

did	name
109	20th Century Fox
110	Bavaria Atelier
101	British Lion
107	Columbia
102	Jean Luc Godard
113	Luso films
104	Mosfilm
103	Paramount
106	Toho
105	United Artists
111	Walt Disney
112	Warner Bros.
108	Westward

The next example shows how to obtain the union of the tables `distributors` and `actors`, restricting the results to those that begin with the letter W in each table. Only distinct rows are wanted, so the key word `ALL` is omitted.

distributors:	actors:
did name	id name
108 Westward	1 Woody Allen
111 Walt Disney	2 Warren Beatty
112 Warner Bros.	3 Walter Matthau
...	...

```
SELECT distributors.name
FROM distributors
WHERE distributors.name LIKE 'W%'
UNION
SELECT actors.name
FROM actors
WHERE actors.name LIKE 'W%';
```

name
Walt Disney
Walter Matthau
Warner Bros.
Warren Beatty
Westward
Woody Allen

This example shows how to use a function in the `FROM` clause, both with and without a column definition list:

```
CREATE FUNCTION distributors(int) RETURNS SETOF distributors AS $$
```

SELECT

```
SELECT * FROM distributors WHERE did = $1;
$$ LANGUAGE SQL;
```

```
SELECT * FROM distributors(111);
 did |    name
-----+-----
 111 | Walt Disney
```

```
CREATE FUNCTION distributors_2(int) RETURNS SETOF record AS $$
    SELECT * FROM distributors WHERE did = $1;
$$ LANGUAGE SQL;
```

```
SELECT * FROM distributors_2(111) AS (f1 int, f2 text);
 f1 |    f2
-----+-----
 111 | Walt Disney
```

Here is an example of a function with an ordinality column added:

```
SELECT * FROM unnest(ARRAY['a','b','c','d','e','f']) WITH ORDINALITY;
unnest | ordinality
-----+-----
 a     |          1
 b     |          2
 c     |          3
 d     |          4
 e     |          5
 f     |          6
(6 rows)
```

This example shows how to use a simple `WITH` clause:

```
WITH t AS (
    SELECT random() as x FROM generate_series(1, 3)
)
SELECT * FROM t
UNION ALL
SELECT * FROM t;
      x
-----
0.534150459803641
0.520092216785997
0.0735620250925422
0.534150459803641
0.520092216785997
0.0735620250925422
```

Notice that the `WITH` query was evaluated only once, so that we got two sets of the same three random values.

This example uses `WITH RECURSIVE` to find all subordinates (direct or indirect) of the employee Mary, and their level of indirectness, from a table that shows only direct subordinates:

```
WITH RECURSIVE employee_recursive(distance, employee_name, manager_name) AS (
    SELECT 1, employee_name, manager_name
    FROM employee
    WHERE manager_name = 'Mary'
    UNION ALL
    SELECT er.distance + 1, e.employee_name, e.manager_name
    FROM employee_recursive er, employee e
    WHERE er.employee_name = e.manager_name
```

```
)  
SELECT distance, employee_name FROM employee_recursive;
```

Notice the typical form of recursive queries: an initial condition, followed by `UNION`, followed by the recursive part of the query. Be sure that the recursive part of the query will eventually return no tuples, or else the query will loop indefinitely. (See [Section 7.8](#) for more examples.)

This example uses `LATERAL` to apply a set-returning function `get_product_names()` for each row of the `manufacturers` table:

```
SELECT m.name AS mname, pname  
FROM manufacturers m, LATERAL get_product_names(m.id) pname;
```

Manufacturers not currently having any products would not appear in the result, since it is an inner join. If we wished to include the names of such manufacturers in the result, we could do:

```
SELECT m.name AS mname, pname  
FROM manufacturers m LEFT JOIN LATERAL get_product_names(m.id) pname ON true;
```

Compatibility

Of course, the `SELECT` statement is compatible with the SQL standard. But there are some extensions and some missing features.

Omitted `FROM` Clauses

Postgres Pro allows one to omit the `FROM` clause. It has a straightforward use to compute the results of simple expressions:

```
SELECT 2+2;
```

```
?column?  
-----  
4
```

Some other SQL databases cannot do this except by introducing a dummy one-row table from which to do the `SELECT`.

Empty `SELECT` Lists

The list of output expressions after `SELECT` can be empty, producing a zero-column result table. This is not valid syntax according to the SQL standard. Postgres Pro allows it to be consistent with allowing zero-column tables. However, an empty list is not allowed when `DISTINCT` is used.

Omitting the `AS` Key Word

In the SQL standard, the optional key word `AS` can be omitted before an output column name whenever the new column name is a valid column name (that is, not the same as any reserved keyword). Postgres Pro is slightly more restrictive: `AS` is required if the new column name matches any keyword at all, reserved or not. Recommended practice is to use `AS` or double-quote output column names, to prevent any possible conflict against future keyword additions.

In `FROM` items, both the standard and Postgres Pro allow `AS` to be omitted before an alias that is an unreserved keyword. But this is impractical for output column names, because of syntactic ambiguities.

Omitting Sub-`SELECT` Aliases in `FROM`

According to the SQL standard, a sub-`SELECT` in the `FROM` list must have an alias. In Postgres Pro, this alias may be omitted.

`ONLY` and Inheritance

The SQL standard requires parentheses around the table name when writing `ONLY`, for example `SELECT * FROM ONLY (tab1), ONLY (tab2) WHERE ...`. Postgres Pro considers these parentheses to be optional.

Postgres Pro allows a trailing `*` to be written to explicitly specify the non-ONLY behavior of including child tables. The standard does not allow this.

(These points apply equally to all SQL commands supporting the ONLY option.)

`TABLESAMPLE` **Clause Restrictions**

The `TABLESAMPLE` clause is currently accepted only on regular tables and materialized views. According to the SQL standard it should be possible to apply it to any `FROM` item.

Function Calls in `FROM`

Postgres Pro allows a function call to be written directly as a member of the `FROM` list. In the SQL standard it would be necessary to wrap such a function call in a sub-`SELECT`; that is, the syntax `FROM func(...) alias` is approximately equivalent to `FROM LATERAL (SELECT func(...)) alias`. Note that `LATERAL` is considered to be implicit; this is because the standard requires `LATERAL` semantics for an `UNNEST()` item in `FROM`. Postgres Pro treats `UNNEST()` the same as other set-returning functions.

Namespace Available to `GROUP BY` **and** `ORDER BY`

In the SQL-92 standard, an `ORDER BY` clause can only use output column names or numbers, while a `GROUP BY` clause can only use expressions based on input column names. Postgres Pro extends each of these clauses to allow the other choice as well (but it uses the standard's interpretation if there is ambiguity). Postgres Pro also allows both clauses to specify arbitrary expressions. Note that names appearing in an expression will always be taken as input-column names, not as output-column names.

SQL:1999 and later use a slightly different definition which is not entirely upward compatible with SQL-92. In most cases, however, Postgres Pro will interpret an `ORDER BY` or `GROUP BY` expression the same way SQL:1999 does.

Functional Dependencies

Postgres Pro recognizes functional dependency (allowing columns to be omitted from `GROUP BY`) only when a table's primary key is included in the `GROUP BY` list. The SQL standard specifies additional conditions that should be recognized.

`LIMIT` **and** `OFFSET`

The clauses `LIMIT` and `OFFSET` are Postgres Pro-specific syntax, also used by MySQL. The SQL:2008 standard has introduced the clauses `OFFSET ... FETCH {FIRST|NEXT} ...` for the same functionality, as shown above in [LIMIT Clause](#). This syntax is also used by IBM DB2. (Applications written for Oracle frequently use a workaround involving the automatically generated `rownum` column, which is not available in Postgres Pro, to implement the effects of these clauses.)

`FOR NO KEY UPDATE`, `FOR UPDATE`, `FOR SHARE`, `FOR KEY SHARE`

Although `FOR UPDATE` appears in the SQL standard, the standard allows it only as an option of `DECLARE CURSOR`. Postgres Pro allows it in any `SELECT` query as well as in sub-`SELECT`s, but this is an extension. The `FOR NO KEY UPDATE`, `FOR SHARE` and `FOR KEY SHARE` variants, as well as the `NOWAIT` and `SKIP LOCKED` options, do not appear in the standard.

Data-Modifying Statements in `WITH`

Postgres Pro allows `INSERT`, `UPDATE`, `DELETE`, and `MERGE` to be used as `WITH` queries. This is not found in the SQL standard.

Nonstandard Clauses

`DISTINCT ON (...)` is an extension of the SQL standard.

`ROWS FROM(...)` is an extension of the SQL standard.

The `MATERIALIZED` and `NOT MATERIALIZED` options of `WITH` are extensions of the SQL standard.

SELECT INTO

SELECT INTO — define a new table from the results of a query

Synopsis

```
[ WITH [ RECURSIVE ] with_query [, ...] ]
SELECT [ ALL | DISTINCT [ ON ( expression [, ...] ) ] ]
    [ { * | expression [ [ AS ] output_name ] } [, ...] ]
    INTO [ TEMPORARY | TEMP | UNLOGGED | CONSTANT ] [ TABLE ] new_table
    [ FROM from_item [, ...] ]
    [ WHERE condition ]
    [ GROUP BY expression [, ...] ]
    [ HAVING condition ]
    [ WINDOW window_name AS ( window_definition ) [, ...] ]
    [ { UNION | INTERSECT | EXCEPT } [ ALL | DISTINCT ] select ]
    [ ORDER BY expression [ ASC | DESC | USING operator ] [ NULLS { FIRST | LAST } ]
    [, ...] ]
    [ LIMIT { count | ALL } ]
    [ OFFSET start [ ROW | ROWS ] ]
    [ FETCH { FIRST | NEXT } [ count ] { ROW | ROWS } ONLY ]
    [ FOR { UPDATE | SHARE } [ OF table_name [, ...] ] [ NOWAIT ] [...] ]
```

Description

SELECT INTO creates a new table and fills it with data computed by a query. The data is not returned to the client, as it is with a normal SELECT. The new table's columns have the names and data types associated with the output columns of the SELECT.

Parameters

TEMPORARY or TEMP

If specified, the table is created as a temporary table. Refer to [CREATE TABLE](#) for details.

UNLOGGED

If specified, the table is created as an unlogged table. Refer to [CREATE TABLE](#) for details.

CONSTANT

If specified, the table is created as read-only. Refer to [CREATE TABLE](#) for details.

new_table

The name (optionally schema-qualified) of the table to be created.

All other parameters are described in detail under [SELECT](#).

Notes

CREATE TABLE AS is functionally similar to SELECT INTO. CREATE TABLE AS is the recommended syntax, since this form of SELECT INTO is not available in ECPG or PL/pgSQL, because they interpret the INTO clause differently. Furthermore, CREATE TABLE AS offers a superset of the functionality provided by SELECT INTO.

In contrast to CREATE TABLE AS, SELECT INTO does not allow specifying properties like a table's access method with USING *method* or the table's tablespace with TABLESPACE *tablespace_name*. Use CREATE TABLE AS if necessary. Therefore, the default table access method is chosen for the new table. See [default_table_access_method](#) for more information.

Examples

Create a new table `films_recent` consisting of only recent entries from the table `films`:

```
SELECT * INTO films_recent FROM films WHERE date_prod >= '2002-01-01';
```

Compatibility

The SQL standard uses `SELECT INTO` to represent selecting values into scalar variables of a host program, rather than creating a new table. This indeed is the usage found in ECPG (see [Chapter 37](#)) and PL/pgSQL (see [Chapter 44](#)). The Postgres Pro usage of `SELECT INTO` to represent table creation is historical. Some other SQL implementations also use `SELECT INTO` in this way (but most SQL implementations support `CREATE TABLE AS` instead). Apart from such compatibility considerations, it is best to use `CREATE TABLE AS` for this purpose in new code.

See Also

[CREATE TABLE AS](#)

SET

SET — change a run-time parameter

Synopsis

```
SET [ SESSION | LOCAL ] configuration_parameter { TO | = } { value | 'value' |  
  DEFAULT }  
SET [ SESSION | LOCAL ] TIME ZONE { value | 'value' | LOCAL | DEFAULT }
```

Description

The `SET` command changes run-time configuration parameters. Many of the run-time parameters listed in [Chapter 19](#) can be changed on-the-fly with `SET`. (Some parameters can only be changed by superusers and users who have been granted `SET` privilege on that parameter. There are also parameters that cannot be changed after server or session start.) `SET` only affects the value used by the current session.

If `SET` (or equivalently `SET SESSION`) is issued within a transaction that is later aborted, the effects of the `SET` command disappear when the transaction is rolled back. Once the surrounding transaction is committed, the effects will persist until the end of the session, unless overridden by another `SET`.

The effects of `SET LOCAL` last only till the end of the current transaction, whether committed or not. A special case is `SET` followed by `SET LOCAL` within a single transaction: the `SET LOCAL` value will be seen until the end of the transaction, but afterwards (if the transaction is committed) the `SET` value will take effect.

The effects of `SET` or `SET LOCAL` are also canceled by rolling back to a savepoint that is earlier than the command.

If `SET LOCAL` is used within a function that has a `SET` option for the same variable (see [CREATE FUNCTION](#)), the effects of the `SET LOCAL` command disappear at function exit; that is, the value in effect when the function was called is restored anyway. This allows `SET LOCAL` to be used for dynamic or repeated changes of a parameter within a function, while still having the convenience of using the `SET` option to save and restore the caller's value. However, a regular `SET` command overrides any surrounding function's `SET` option; its effects will persist unless rolled back.

Note

In PostgreSQL versions 8.0 through 8.2, the effects of a `SET LOCAL` would be canceled by releasing an earlier savepoint, or by successful exit from a PL/pgSQL exception block. This behavior has been changed because it was deemed unintuitive.

Parameters

`SESSION`

Specifies that the command takes effect for the current session. (This is the default if neither `SESSION` nor `LOCAL` appears.)

`LOCAL`

Specifies that the command takes effect for only the current transaction. After `COMMIT` or `ROLLBACK`, the session-level setting takes effect again. Issuing this outside of a transaction block emits a warning and otherwise has no effect.

configuration_parameter

Name of a settable run-time parameter. Available parameters are documented in [Chapter 19](#) and below.

value

New value of parameter. Values can be specified as string constants, identifiers, numbers, or comma-separated lists of these, as appropriate for the particular parameter. `DEFAULT` can be written to specify resetting the parameter to its default value (that is, whatever value it would have had if no `SET` had been executed in the current session).

Besides the configuration parameters documented in [Chapter 19](#), there are a few that can only be adjusted using the `SET` command or that have a special syntax:

SCHEMA

`SET SCHEMA 'value'` is an alias for `SET search_path TO value`. Only one schema can be specified using this syntax.

NAMES

`SET NAMES value` is an alias for `SET client_encoding TO value`.

SEED

Sets the internal seed for the random number generator (the function `random`). Allowed values are floating-point numbers between -1 and 1 inclusive.

The seed can also be set by invoking the function `setseed`:

```
SELECT setseed(value);
```

TIME ZONE

`SET TIME ZONE 'value'` is an alias for `SET timezone TO 'value'`. The syntax `SET TIME ZONE` allows special syntax for the time zone specification. Here are examples of valid values:

`'America/Los_Angeles'`

The time zone for Berkeley, California.

`'Europe/Rome'`

The time zone for Italy.

`-7`

The time zone 7 hours west from UTC (equivalent to PDT). Positive values are east from UTC.

`INTERVAL '-08:00' HOUR TO MINUTE`

The time zone 8 hours west from UTC (equivalent to PST).

LOCAL

DEFAULT

Set the time zone to your local time zone (that is, the server's default value of `timezone`).

Timezone settings given as numbers or intervals are internally translated to POSIX timezone syntax. For example, after `SET TIME ZONE -7`, `SHOW TIME ZONE` would report `<-07>+07`.

Time zone abbreviations are not supported by `SET`; see [Section 8.5.3](#) for more information about time zones.

Notes

The function `set_config` provides equivalent functionality; see [Section 9.28.1](#). Also, it is possible to `UPDATE` the `pg_settings` system view to perform the equivalent of `SET`.

Examples

Set the schema search path:

```
SET search_path TO my_schema, public;
```

Set the style of date to traditional POSTGRES with “day before month” input convention:

```
SET datestyle TO postgres, dmy;
```

Set the time zone for Berkeley, California:

```
SET TIME ZONE 'America/Los_Angeles';
```

Set the time zone for Italy:

```
SET TIME ZONE 'Europe/Rome';
```

Compatibility

`SET TIME ZONE` extends syntax defined in the SQL standard. The standard allows only numeric time zone offsets while Postgres Pro allows more flexible time-zone specifications. All other `SET` features are Postgres Pro extensions.

See Also

[RESET](#), [SHOW](#)

SET CONSTRAINTS

SET CONSTRAINTS — set constraint check timing for the current transaction

Synopsis

```
SET CONSTRAINTS { ALL | name [, ...] } { DEFERRED | IMMEDIATE }
```

Description

SET CONSTRAINTS sets the behavior of constraint checking within the current transaction. IMMEDIATE constraints are checked at the end of each statement. DEFERRED constraints are not checked until transaction commit. Each constraint has its own IMMEDIATE or DEFERRED mode.

Upon creation, a constraint is given one of three characteristics: DEFERRABLE INITIALLY DEFERRED, DEFERRABLE INITIALLY IMMEDIATE, or NOT DEFERRABLE. The third class is always IMMEDIATE and is not affected by the SET CONSTRAINTS command. The first two classes start every transaction in the indicated mode, but their behavior can be changed within a transaction by SET CONSTRAINTS.

SET CONSTRAINTS with a list of constraint names changes the mode of just those constraints (which must all be deferrable). Each constraint name can be schema-qualified. The current schema search path is used to find the first matching name if no schema name is specified. SET CONSTRAINTS ALL changes the mode of all deferrable constraints.

When SET CONSTRAINTS changes the mode of a constraint from DEFERRED to IMMEDIATE, the new mode takes effect retroactively: any outstanding data modifications that would have been checked at the end of the transaction are instead checked during the execution of the SET CONSTRAINTS command. If any such constraint is violated, the SET CONSTRAINTS fails (and does not change the constraint mode). Thus, SET CONSTRAINTS can be used to force checking of constraints to occur at a specific point in a transaction.

Currently, only UNIQUE, PRIMARY KEY, REFERENCES (foreign key), and EXCLUDE constraints are affected by this setting. NOT NULL and CHECK constraints are always checked immediately when a row is inserted or modified (*not* at the end of the statement). Uniqueness and exclusion constraints that have not been declared DEFERRABLE are also checked immediately.

The firing of triggers that are declared as “constraint triggers” is also controlled by this setting — they fire at the same time that the associated constraint should be checked.

Notes

Because Postgres Pro does not require constraint names to be unique within a schema (but only per-table), it is possible that there is more than one match for a specified constraint name. In this case SET CONSTRAINTS will act on all matches. For a non-schema-qualified name, once a match or matches have been found in some schema in the search path, schemas appearing later in the path are not searched.

This command only alters the behavior of constraints within the current transaction. Issuing this outside of a transaction block emits a warning and otherwise has no effect.

Compatibility

This command complies with the behavior defined in the SQL standard, except for the limitation that, in Postgres Pro, it does not apply to NOT NULL and CHECK constraints. Also, Postgres Pro checks non-deferrable uniqueness constraints immediately, not at end of statement as the standard would suggest.

SET ROLE

SET ROLE — set the current user identifier of the current session

Synopsis

```
SET [ SESSION | LOCAL ] ROLE role_name
SET [ SESSION | LOCAL ] ROLE NONE
RESET ROLE
```

Description

This command sets the current user identifier of the current SQL session to be *role_name*. The role name can be written as either an identifier or a string literal. After SET ROLE, permissions checking for SQL commands is carried out as though the named role were the one that had logged in originally. Note that SET ROLE and SET SESSION AUTHORIZATION are exceptions; permissions checks for those continue to use the current session user and the initial session user (the *authenticated user*), respectively.

The current session user must have the SET option for the specified *role_name*, either directly or indirectly via a chain of memberships with the SET option. (If the session user is a superuser, any role can be selected.)

The SESSION and LOCAL modifiers act the same as for the regular SET command.

SET ROLE NONE sets the current user identifier to the current session user identifier, as returned by session_user. RESET ROLE sets the current user identifier to the connection-time setting specified by the [command-line options](#), ALTER ROLE, or ALTER DATABASE, if any such settings exist. Otherwise, RESET ROLE sets the current user identifier to the current session user identifier. These forms can be executed by any user.

Notes

Using this command, it is possible to either add privileges or restrict one's privileges. If the session user role has been granted memberships WITH INHERIT TRUE, it automatically has all the privileges of every such role. In this case, SET ROLE effectively drops all the privileges except for those which the target role directly possesses or inherits. On the other hand, if the session user role has been granted memberships WITH INHERIT FALSE, the privileges of the granted roles can't be accessed by default. However, if the role was granted WITH SET TRUE, the session user can use SET ROLE to drop the privileges assigned directly to the session user and instead acquire the privileges available to the named role. If the role was granted WITH INHERIT FALSE, SET FALSE then the privileges of that role cannot be exercised either with or without SET ROLE.

SET ROLE has effects comparable to SET SESSION AUTHORIZATION, but the privilege checks involved are quite different. Also, SET SESSION AUTHORIZATION determines which roles are allowable for later SET ROLE commands, whereas changing roles with SET ROLE does not change the set of roles allowed to a later SET ROLE.

SET ROLE does not process session variables as specified by the role's ALTER ROLE settings; this only happens during login.

SET ROLE cannot be used within a SECURITY DEFINER function.

Examples

```
SELECT SESSION_USER, CURRENT_USER;
```

```
session_user | current_user
-----+-----
```



```
peter          | peter

SET ROLE 'paul';

SELECT SESSION_USER, CURRENT_USER;

 session_user | current_user
-----+-----
peter        | paul
```

Compatibility

Postgres Pro allows identifier syntax ("*rolename*"), while the SQL standard requires the role name to be written as a string literal. SQL does not allow this command during a transaction; Postgres Pro does not make this restriction because there is no reason to. The `SESSION` and `LOCAL` modifiers are a Postgres Pro extension, as is the `RESET` syntax.

See Also

[SET SESSION AUTHORIZATION](#)

SET SESSION AUTHORIZATION

SET SESSION AUTHORIZATION — set the session user identifier and the current user identifier of the current session

Synopsis

```
SET [ SESSION | LOCAL ] SESSION AUTHORIZATION user_name
SET [ SESSION | LOCAL ] SESSION AUTHORIZATION DEFAULT
RESET SESSION AUTHORIZATION
```

Description

This command sets the session user identifier and the current user identifier of the current SQL session to be *user_name*. The user name can be written as either an identifier or a string literal. Using this command, it is possible, for example, to temporarily become an unprivileged user and later switch back to being a superuser.

The session user identifier is initially set to be the (possibly authenticated) user name provided by the client. The current user identifier is normally equal to the session user identifier, but might change temporarily in the context of SECURITY DEFINER functions and similar mechanisms; it can also be changed by SET ROLE. The current user identifier is relevant for permission checking.

The session user identifier can be changed only if the initial session user (the *authenticated user*) has the superuser privilege. Otherwise, the command is accepted only if it specifies the authenticated user name.

The SESSION and LOCAL modifiers act the same as for the regular SET command.

The DEFAULT and RESET forms reset the session and current user identifiers to be the originally authenticated user name. These forms can be executed by any user.

Notes

SET SESSION AUTHORIZATION cannot be used within a SECURITY DEFINER function.

Examples

```
SELECT SESSION_USER, CURRENT_USER;
```

```
 session_user | current_user
-----+-----
peter        | peter
```

```
SET SESSION AUTHORIZATION 'paul';
```

```
SELECT SESSION_USER, CURRENT_USER;
```

```
 session_user | current_user
-----+-----
paul         | paul
```

Compatibility

The SQL standard allows some other expressions to appear in place of the literal *user_name*, but these options are not important in practice. Postgres Pro allows identifier syntax ("*username*"), which SQL does not. SQL does not allow this command during a transaction; Postgres Pro does not make this restriction because there is no reason to. The SESSION and LOCAL modifiers are a Postgres Pro extension, as is the RESET syntax.

The privileges necessary to execute this command are left implementation-defined by the standard.

See Also

[SET ROLE](#)

SET TRANSACTION

SET TRANSACTION — set the characteristics of the current transaction

Synopsis

```
SET TRANSACTION transaction_mode [, ...]
SET TRANSACTION SNAPSHOT snapshot_id
SET SESSION CHARACTERISTICS AS TRANSACTION transaction_mode [, ...]
```

where *transaction_mode* is one of:

```
ISOLATION LEVEL { SERIALIZABLE | REPEATABLE READ | READ COMMITTED | READ
UNCOMMITTED }
READ WRITE | READ ONLY
[ NOT ] DEFERRABLE
```

Description

The `SET TRANSACTION` command sets the characteristics of the current transaction. It has no effect on any subsequent transactions. `SET SESSION CHARACTERISTICS` sets the default transaction characteristics for subsequent transactions of a session. These defaults can be overridden by `SET TRANSACTION` for an individual transaction.

The available transaction characteristics are the transaction isolation level, the transaction access mode (read/write or read-only), and the deferrable mode. In addition, a snapshot can be selected, though only for the current transaction, not as a session default.

The isolation level of a transaction determines what data the transaction can see when other transactions are running concurrently:

`READ COMMITTED`

A statement can only see rows committed before it began. This is the default.

`REPEATABLE READ`

All statements of the current transaction can only see rows committed before the first query or data-modification statement was executed in this transaction.

`SERIALIZABLE`

All statements of the current transaction can only see rows committed before the first query or data-modification statement was executed in this transaction. If a pattern of reads and writes among concurrent serializable transactions would create a situation which could not have occurred for any serial (one-at-a-time) execution of those transactions, one of them will be rolled back with a `serialization_failure` error.

The SQL standard defines one additional level, `READ UNCOMMITTED`. In Postgres Pro `READ UNCOMMITTED` is treated as `READ COMMITTED`.

The transaction isolation level cannot be changed after the first query or data-modification statement (`SELECT`, `INSERT`, `DELETE`, `UPDATE`, `MERGE`, `FETCH`, or `COPY`) of a transaction has been executed. See [Chapter 13](#) for more information about transaction isolation and concurrency control.

The transaction access mode determines whether the transaction is read/write or read-only. Read/write is the default. When a transaction is read-only, the following SQL commands are disallowed: `INSERT`, `UPDATE`, `DELETE`, `MERGE`, and `COPY FROM` if the table they would write to is not a temporary table; all `CREATE`, `ALTER`, and `DROP` commands; `COMMENT`, `GRANT`, `REVOKE`, `TRUNCATE`; and `EXPLAIN ANALYZE` and `EXECUTE` if the command they would execute is among those listed. This is a high-level notion of read-only that does not prevent all writes to disk.

The `DEFERRABLE` transaction property has no effect unless the transaction is also `SERIALIZABLE` and `READ ONLY`. When all three of these properties are selected for a transaction, the transaction may block when first acquiring its snapshot, after which it is able to run without the normal overhead of a `SERIALIZABLE` transaction and without any risk of contributing to or being canceled by a serialization failure. This mode is well suited for long-running reports or backups.

The `SET TRANSACTION SNAPSHOT` command allows a new transaction to run with the same *snapshot* as an existing transaction. The pre-existing transaction must have exported its snapshot with the `pg_export_snapshot` function (see [Section 9.28.5](#)). That function returns a snapshot identifier, which must be given to `SET TRANSACTION SNAPSHOT` to specify which snapshot is to be imported. The identifier must be written as a string literal in this command, for example `'00000003-0000001B-1'`. `SET TRANSACTION SNAPSHOT` can only be executed at the start of a transaction, before the first query or data-modification statement (`SELECT`, `INSERT`, `DELETE`, `UPDATE`, `MERGE`, `FETCH`, or `COPY`) of the transaction. Furthermore, the transaction must already be set to `SERIALIZABLE` or `REPEATABLE READ` isolation level (otherwise, the snapshot would be discarded immediately, since `READ COMMITTED` mode takes a new snapshot for each command). If the importing transaction uses `SERIALIZABLE` isolation level, then the transaction that exported the snapshot must also use that isolation level. Also, a non-read-only serializable transaction cannot import a snapshot from a read-only transaction.

Notes

If `SET TRANSACTION` is executed without a prior `START TRANSACTION` or `BEGIN`, it emits a warning and otherwise has no effect.

It is possible to dispense with `SET TRANSACTION` by instead specifying the desired *transaction_modes* in `BEGIN` or `START TRANSACTION`. But that option is not available for `SET TRANSACTION SNAPSHOT`.

The session default transaction modes can also be set or examined via the configuration parameters [default_transaction_isolation](#), [default_transaction_read_only](#), and [default_transaction_deferrable](#). (In fact `SET SESSION CHARACTERISTICS` is just a verbose equivalent for setting these variables with `SET`.) This means the defaults can be set in the configuration file, via `ALTER DATABASE`, etc. Consult [Chapter 19](#) for more information.

The current transaction's modes can similarly be set or examined via the configuration parameters [transaction_isolation](#), [transaction_read_only](#), and [transaction_deferrable](#). Setting one of these parameters acts the same as the corresponding `SET TRANSACTION` option, with the same restrictions on when it can be done. However, these parameters cannot be set in the configuration file, or from any source other than live SQL.

Examples

To begin a new transaction with the same snapshot as an already existing transaction, first export the snapshot from the existing transaction. That will return the snapshot identifier, for example:

```
BEGIN TRANSACTION ISOLATION LEVEL REPEATABLE READ;
SELECT pg_export_snapshot();
pg_export_snapshot
-----
00000003-0000001B-1
(1 row)
```

Then give the snapshot identifier in a `SET TRANSACTION SNAPSHOT` command at the beginning of the newly opened transaction:

```
BEGIN TRANSACTION ISOLATION LEVEL REPEATABLE READ;
SET TRANSACTION SNAPSHOT '00000003-0000001B-1';
```

Compatibility

These commands are defined in the SQL standard, except for the `DEFERRABLE` transaction mode and the `SET TRANSACTION SNAPSHOT` form, which are Postgres Pro extensions.

`SERIALIZABLE` is the default transaction isolation level in the standard. In Postgres Pro the default is ordinarily `READ COMMITTED`, but you can change it as mentioned above.

In the SQL standard, there is one other transaction characteristic that can be set with these commands: the size of the diagnostics area. This concept is specific to embedded SQL, and therefore is not implemented in the Postgres Pro server.

The SQL standard requires commas between successive *transaction_modes*, but for historical reasons Postgres Pro allows the commas to be omitted.

SHOW

SHOW — show the value of a run-time parameter

Synopsis

```
SHOW name
SHOW ALL
```

Description

SHOW will display the current setting of run-time parameters. These variables can be set using the SET statement, by editing the `postgresql.conf` configuration file, through the `PGOPTIONS` environmental variable (when using libpq or a libpq-based application), or through command-line flags when starting the `postgres` server. See [Chapter 19](#) for details.

Parameters

name

The name of a run-time parameter. Available parameters are documented in [Chapter 19](#) and on the [SET](#) reference page. In addition, there are a few parameters that can be shown but not set:

SERVER_VERSION

Shows the server's version number.

SERVER_ENCODING

Shows the server-side character set encoding. At present, this parameter can be shown but not set, because the encoding is determined at database creation time.

IS_SUPERUSER

True if the current role has superuser privileges.

ALL

Show the values of all configuration parameters, with descriptions.

Notes

The function `current_setting` produces equivalent output; see [Section 9.28.1](#). Also, the `pg_settings` system view produces the same information.

Examples

Show the current setting of the parameter `DateStyle`:

```
SHOW DateStyle;
DateStyle
-----
ISO, MDY
(1 row)
```

Show the current setting of the parameter `geqo`:

```
SHOW geqo;
geqo
-----
on
(1 row)
```

Show all settings:

```
SHOW ALL;
```

name	setting	description
allow_system_table_mods	off	Allows modifications of the structure of ...
.	.	.
xmloption	content	Sets whether XML data in implicit parsing ...
zero_damaged_pages	off	Continues processing past damaged page headers.

(196 rows)

Compatibility

The SHOW command is a Postgres Pro extension.

See Also

[SET](#), [RESET](#)

START TRANSACTION

START TRANSACTION — start a transaction block

Synopsis

```
START [ AUTONOMOUS ] TRANSACTION [ transaction_mode [, ...] ]
```

where *transaction_mode* is one of:

```
ISOLATION LEVEL { SERIALIZABLE | REPEATABLE READ | READ COMMITTED | READ  
UNCOMMITTED }  
READ WRITE | READ ONLY  
[ NOT ] DEFERRABLE
```

Description

This command begins a new transaction block. If the isolation level, read/write mode, or deferrable mode is specified, the new transaction has those characteristics, as if [SET TRANSACTION](#) was executed. This is the same as the [BEGIN](#) command.

If you specify an optional `AUTONOMOUS` keyword, an autonomous transaction is started. Autonomous transactions can be started only within another transaction. For details, see [Chapter 16](#).

Parameters

Refer to [SET TRANSACTION](#) for information on the meaning of the parameters to this statement.

Compatibility

In the standard, it is not necessary to issue `START TRANSACTION` to start a transaction block: any SQL command implicitly begins a block. Postgres Pro's behavior can be seen as implicitly issuing a `COMMIT` after each command that does not follow `START TRANSACTION` (or `BEGIN`), and it is therefore often called “autocommit”. Other relational database systems might offer an autocommit feature as a convenience.

The `DEFERRABLE transaction_mode` is a Postgres Pro language extension.

The SQL standard requires commas between successive *transaction_modes*, but for historical reasons Postgres Pro allows the commas to be omitted.

See also the compatibility section of [SET TRANSACTION](#).

See Also

[BEGIN](#), [COMMIT](#), [ROLLBACK](#), [SAVEPOINT](#), [SET TRANSACTION](#)

TRUNCATE

TRUNCATE — empty a table or set of tables

Synopsis

```
TRUNCATE [ TABLE ] [ ONLY ] name [ * ] [, ... ]  
        [ RESTART IDENTITY | CONTINUE IDENTITY ] [ CASCADE | RESTRICT ]
```

Description

TRUNCATE quickly removes all rows from a set of tables. It has the same effect as an unqualified DELETE on each table, but since it does not actually scan the tables it is faster. Furthermore, it reclaims disk space immediately, rather than requiring a subsequent VACUUM operation. This is most useful on large tables.

Parameters

name

The name (optionally schema-qualified) of a table to truncate. If ONLY is specified before the table name, only that table is truncated. If ONLY is not specified, the table and all its descendant tables (if any) are truncated. Optionally, * can be specified after the table name to explicitly indicate that descendant tables are included.

RESTART IDENTITY

Automatically restart sequences owned by columns of the truncated table(s).

CONTINUE IDENTITY

Do not change the values of sequences. This is the default.

CASCADE

Automatically truncate all tables that have foreign-key references to any of the named tables, or to any tables added to the group due to CASCADE.

RESTRICT

Refuse to truncate if any of the tables have foreign-key references from tables that are not listed in the command. This is the default.

Notes

You must have the TRUNCATE privilege on a table to truncate it.

TRUNCATE acquires an ACCESS EXCLUSIVE lock on each table it operates on, which blocks all other concurrent operations on the table. When RESTART IDENTITY is specified, any sequences that are to be restarted are likewise locked exclusively. If concurrent access to a table is required, then the DELETE command should be used instead.

TRUNCATE cannot be used on a table that has foreign-key references from other tables, unless all such tables are also truncated in the same command. Checking validity in such cases would require table scans, and the whole point is not to do one. The CASCADE option can be used to automatically include all dependent tables — but be very careful when using this option, or else you might lose data you did not intend to! Note in particular that when the table to be truncated is a partition, siblings partitions are left untouched, but cascading occurs to all referencing tables and all their partitions with no distinction.

TRUNCATE will not fire any ON DELETE triggers that might exist for the tables. But it will fire ON TRUNCATE triggers. If ON TRUNCATE triggers are defined for any of the tables, then all BEFORE TRUNCATE triggers are

fired before any truncation happens, and all `AFTER TRUNCATE` triggers are fired after the last truncation is performed and any sequences are reset. The triggers will fire in the order that the tables are to be processed (first those listed in the command, and then any that were added due to cascading).

`TRUNCATE` is not MVCC-safe. After truncation, the table will appear empty to concurrent transactions, if they are using a snapshot taken before the truncation occurred. See [Section 13.6](#) for more details.

`TRUNCATE` is transaction-safe with respect to the data in the tables: the truncation will be safely rolled back if the surrounding transaction does not commit.

When `RESTART IDENTITY` is specified, the implied `ALTER SEQUENCE RESTART` operations are also done transactionally; that is, they will be rolled back if the surrounding transaction does not commit. Be aware that if any additional sequence operations are done on the restarted sequences before the transaction rolls back, the effects of these operations on the sequences will be rolled back, but not their effects on `currval()`; that is, after the transaction `currval()` will continue to reflect the last sequence value obtained inside the failed transaction, even though the sequence itself may no longer be consistent with that. This is similar to the usual behavior of `currval()` after a failed transaction.

`TRUNCATE` can be used for foreign tables if supported by the foreign data wrapper, for instance, see [postgres_fdw](#).

Examples

Truncate the tables `bigtable` and `fattable`:

```
TRUNCATE bigtable, fattable;
```

The same, and also reset any associated sequence generators:

```
TRUNCATE bigtable, fattable RESTART IDENTITY;
```

Truncate the table `othertable`, and cascade to any tables that reference `othertable` via foreign-key constraints:

```
TRUNCATE othertable CASCADE;
```

Compatibility

The SQL:2008 standard includes a `TRUNCATE` command with the syntax `TRUNCATE TABLE tablename`. The clauses `CONTINUE IDENTITY/RESTART IDENTITY` also appear in that standard, but have slightly different though related meanings. Some of the concurrency behavior of this command is left implementation-defined by the standard, so the above notes should be considered and compared with other implementations if necessary.

See Also

[DELETE](#)

UNLISTEN

UNLISTEN — stop listening for a notification

Synopsis

```
UNLISTEN { channel | * }
```

Description

UNLISTEN is used to remove an existing registration for NOTIFY events. UNLISTEN cancels any existing registration of the current Postgres Pro session as a listener on the notification channel named *channel*. The special wildcard * cancels all listener registrations for the current session.

[NOTIFY](#) contains a more extensive discussion of the use of LISTEN and NOTIFY.

Parameters

channel

Name of a notification channel (any identifier).

*

All current listen registrations for this session are cleared.

Notes

You can unlisten something you were not listening for; no warning or error will appear.

At the end of each session, UNLISTEN * is automatically executed.

A transaction that has executed UNLISTEN cannot be prepared for two-phase commit.

Examples

To make a registration:

```
LISTEN virtual;
NOTIFY virtual;
Asynchronous notification "virtual" received from server process with PID 8448.
```

Once UNLISTEN has been executed, further NOTIFY messages will be ignored:

```
UNLISTEN virtual;
NOTIFY virtual;
-- no NOTIFY event is received
```

Compatibility

There is no UNLISTEN command in the SQL standard.

See Also

[LISTEN](#), [NOTIFY](#)

UPDATE

UPDATE — update rows of a table

Synopsis

```
[ WITH [ RECURSIVE ] with_query [, ...] ]
UPDATE [ ONLY ] table_name [ * ] [ [ AS ] alias ]
    SET { column_name = { expression | DEFAULT } |
        ( column_name [, ...] ) = [ ROW ] ( { expression | DEFAULT } [, ...] ) |
        ( column_name [, ...] ) = ( sub-SELECT )
    } [, ...]
[ FROM from_item [, ...] ]
[ WHERE condition | WHERE CURRENT OF cursor_name ]
[ RETURNING { * | output_expression [ [ AS ] output_name ] } [, ...] ]
```

Description

UPDATE changes the values of the specified columns in all rows that satisfy the condition. Only the columns to be modified need be mentioned in the SET clause; columns not explicitly modified retain their previous values.

There are two ways to modify a table using information contained in other tables in the database: using sub-selects, or specifying additional tables in the FROM clause. Which technique is more appropriate depends on the specific circumstances.

The optional RETURNING clause causes UPDATE to compute and return value(s) based on each row actually updated. Any expression using the table's columns, and/or columns of other tables mentioned in FROM, can be computed. The new (post-update) values of the table's columns are used. The syntax of the RETURNING list is identical to that of the output list of SELECT.

You must have the UPDATE privilege on the table, or at least on the column(s) that are listed to be updated. You must also have the SELECT privilege on any column whose values are read in the *expressions* or *condition*.

Parameters

with_query

The WITH clause allows you to specify one or more subqueries that can be referenced by name in the UPDATE query. See [Section 7.8](#) and [SELECT](#) for details.

table_name

The name (optionally schema-qualified) of the table to update. If ONLY is specified before the table name, matching rows are updated in the named table only. If ONLY is not specified, matching rows are also updated in any tables inheriting from the named table. Optionally, * can be specified after the table name to explicitly indicate that descendant tables are included.

alias

A substitute name for the target table. When an alias is provided, it completely hides the actual name of the table. For example, given UPDATE foo AS f, the remainder of the UPDATE statement must refer to this table as f not foo.

column_name

The name of a column in the table named by *table_name*. The column name can be qualified with a subfield name or array subscript, if needed. Do not include the table's name in the specification of a target column — for example, UPDATE table_name SET table_name.col = 1 is invalid.

expression

An expression to assign to the column. The expression can use the old values of this and other columns in the table.

DEFAULT

Set the column to its default value (which will be NULL if no specific default expression has been assigned to it). An identity column will be set to a new value generated by the associated sequence. For a generated column, specifying this is permitted but merely specifies the normal behavior of computing the column from its generation expression.

sub-SELECT

A `SELECT` sub-query that produces as many output columns as are listed in the parenthesized column list preceding it. The sub-query must yield no more than one row when executed. If it yields one row, its column values are assigned to the target columns; if it yields no rows, NULL values are assigned to the target columns. The sub-query can refer to old values of the current row of the table being updated.

from_item

A table expression allowing columns from other tables to appear in the `WHERE` condition and update expressions. This uses the same syntax as the `FROM` clause of a `SELECT` statement; for example, an alias for the table name can be specified. Do not repeat the target table as a *from_item* unless you intend a self-join (in which case it must appear with an alias in the *from_item*).

condition

An expression that returns a value of type `boolean`. Only rows for which this expression returns `true` will be updated.

cursor_name

The name of the cursor to use in a `WHERE CURRENT OF` condition. The row to be updated is the one most recently fetched from this cursor. The cursor must be a non-grouping query on the `UPDATE`'s target table. Note that `WHERE CURRENT OF` cannot be specified together with a Boolean condition. See [DECLARE](#) for more information about using cursors with `WHERE CURRENT OF`.

output_expression

An expression to be computed and returned by the `UPDATE` command after each row is updated. The expression can use any column names of the table named by *table_name* or table(s) listed in `FROM`. Write `*` to return all columns.

output_name

A name to use for a returned column.

Outputs

On successful completion, an `UPDATE` command returns a command tag of the form

```
UPDATE count
```

The *count* is the number of rows updated, including matched rows whose values did not change. Note that the number may be less than the number of rows that matched the *condition* when updates were suppressed by a `BEFORE UPDATE` trigger. If *count* is 0, no rows were updated by the query (this is not considered an error).

If the `UPDATE` command contains a `RETURNING` clause, the result will be similar to that of a `SELECT` statement containing the columns and values defined in the `RETURNING` list, computed over the row(s) updated by the command.

Notes

When a `FROM` clause is present, what essentially happens is that the target table is joined to the tables mentioned in the *from_item* list, and each output row of the join represents an update operation for the target table. When using `FROM` you should ensure that the join produces at most one output row for each row to be modified. In other words, a target row shouldn't join to more than one row from the other table(s). If it does, then only one of the join rows will be used to update the target row, but which one will be used is not readily predictable.

Because of this indeterminacy, referencing other tables only within sub-selects is safer, though often harder to read and slower than using a join.

In the case of a partitioned table, updating a row might cause it to no longer satisfy the partition constraint of the containing partition. In that case, if there is some other partition in the partition tree for which this row satisfies its partition constraint, then the row is moved to that partition. If there is no such partition, an error will occur. Behind the scenes, the row movement is actually a `DELETE` and `INSERT` operation.

There is a possibility that a concurrent `UPDATE` or `DELETE` on the row being moved will get a serialization failure error. Suppose session 1 is performing an `UPDATE` on a partition key, and meanwhile a concurrent session 2 for which this row is visible performs an `UPDATE` or `DELETE` operation on this row. In such case, session 2's `UPDATE` or `DELETE` will detect the row movement and raise a serialization failure error (which always returns with an `SQLSTATE` code '40001'). Applications may wish to retry the transaction if this occurs. In the usual case where the table is not partitioned, or where there is no row movement, session 2 would have identified the newly updated row and carried out the `UPDATE/DELETE` on this new row version.

Note that while rows can be moved from local partitions to a foreign-table partition (provided the foreign data wrapper supports tuple routing), they cannot be moved from a foreign-table partition to another partition.

An attempt of moving a row from one partition to another will fail if a foreign key is found to directly reference an ancestor of the source partition that is not the same as the ancestor that's mentioned in the `UPDATE` query.

Examples

Change the word `Drama` to `Dramatic` in the column `kind` of the table `films`:

```
UPDATE films SET kind = 'Dramatic' WHERE kind = 'Drama';
```

Adjust temperature entries and reset precipitation to its default value in one row of the table `weather`:

```
UPDATE weather SET temp_lo = temp_lo+1, temp_hi = temp_lo+15, prcp = DEFAULT
WHERE city = 'San Francisco' AND date = '2003-07-03';
```

Perform the same operation and return the updated entries:

```
UPDATE weather SET temp_lo = temp_lo+1, temp_hi = temp_lo+15, prcp = DEFAULT
WHERE city = 'San Francisco' AND date = '2003-07-03'
RETURNING temp_lo, temp_hi, prcp;
```

Use the alternative column-list syntax to do the same update:

```
UPDATE weather SET (temp_lo, temp_hi, prcp) = (temp_lo+1, temp_lo+15, DEFAULT)
WHERE city = 'San Francisco' AND date = '2003-07-03';
```

Increment the sales count of the salesperson who manages the account for Acme Corporation, using the `FROM` clause syntax:

```
UPDATE employees SET sales_count = sales_count + 1 FROM accounts
WHERE accounts.name = 'Acme Corporation'
AND employees.id = accounts.sales_person;
```

Perform the same operation, using a sub-select in the `WHERE` clause:

```
UPDATE employees SET sales_count = sales_count + 1 WHERE id =
  (SELECT sales_person FROM accounts WHERE name = 'Acme Corporation');
```

Update contact names in an accounts table to match the currently assigned salespeople:

```
UPDATE accounts SET (contact_first_name, contact_last_name) =
  (SELECT first_name, last_name FROM employees
   WHERE employees.id = accounts.sales_person);
```

A similar result could be accomplished with a join:

```
UPDATE accounts SET contact_first_name = first_name,
                    contact_last_name = last_name
  FROM employees WHERE employees.id = accounts.sales_person;
```

However, the second query may give unexpected results if `employees.id` is not a unique key, whereas the first query is guaranteed to raise an error if there are multiple `id` matches. Also, if there is no match for a particular `accounts.sales_person` entry, the first query will set the corresponding name fields to `NULL`, whereas the second query will not update that row at all.

Update statistics in a summary table to match the current data:

```
UPDATE summary s SET (sum_x, sum_y, avg_x, avg_y) =
  (SELECT sum(x), sum(y), avg(x), avg(y) FROM data d
   WHERE d.group_id = s.group_id);
```

Attempt to insert a new stock item along with the quantity of stock. If the item already exists, instead update the stock count of the existing item. To do this without failing the entire transaction, use savepoints:

```
BEGIN;
-- other operations
SAVEPOINT sp1;
INSERT INTO wines VALUES('Chateau Lafite 2003', '24');
-- Assume the above fails because of a unique key violation,
-- so now we issue these commands:
ROLLBACK TO sp1;
UPDATE wines SET stock = stock + 24 WHERE winename = 'Chateau Lafite 2003';
-- continue with other operations, and eventually
COMMIT;
```

Change the `kind` column of the table `films` in the row on which the cursor `c_films` is currently positioned:

```
UPDATE films SET kind = 'Dramatic' WHERE CURRENT OF c_films;
```

Updates affecting many rows can have negative effects on system performance, such as table bloat, increased replica lag, and increased lock contention. In such situations it can make sense to perform the operation in smaller batches, possibly with a `VACUUM` operation on the table between batches. While there is no `LIMIT` clause for `UPDATE`, it is possible to get a similar effect through the use of a [Common Table Expression](#) and a self-join. With the standard Postgres Pro table access method, a self-join on the system column `ctid` is very efficient:

```
WITH exceeded_max_retries AS (
  SELECT w.ctid FROM work_item AS w
    WHERE w.status = 'active' AND w.num_retries > 10
    ORDER BY w.retry_timestamp
    FOR UPDATE
    LIMIT 5000
)
UPDATE work_item SET status = 'failed'
  FROM exceeded_max_retries AS emr
```



```
WHERE work_item.ctid = emr.ctid;
```

This command will need to be repeated until no rows remain to be updated. Use of an `ORDER BY` clause allows the command to prioritize which rows will be updated; it can also prevent deadlock with other update operations if they use the same ordering. If lock contention is a concern, then `SKIP LOCKED` can be added to the CTE to prevent multiple commands from updating the same row. However, then a final `UPDATE` without `SKIP LOCKED` or `LIMIT` will be needed to ensure that no matching rows were overlooked.

Compatibility

This command conforms to the SQL standard, except that the `FROM` and `RETURNING` clauses are Postgres Pro extensions, as is the ability to use `WITH` with `UPDATE`.

Some other database systems offer a `FROM` option in which the target table is supposed to be listed again within `FROM`. That is not how Postgres Pro interprets `FROM`. Be careful when porting applications that use this extension.

According to the standard, the source value for a parenthesized sub-list of target column names can be any row-valued expression yielding the correct number of columns. Postgres Pro only allows the source value to be a [row constructor](#) or a sub-`SELECT`. An individual column's updated value can be specified as `DEFAULT` in the row-constructor case, but not inside a sub-`SELECT`.

VACUUM

VACUUM — garbage-collect and optionally analyze a database

Synopsis

```
VACUUM [ ( option [, ...] ) ] [ table_and_columns [, ...] ]
```

where *option* can be one of:

```
FULL [ boolean ]
FREEZE [ boolean ]
VERBOSE [ boolean ]
ANALYZE [ boolean ]
DISABLE_PAGE_SKIPPING [ boolean ]
SKIP_LOCKED [ boolean ]
INDEX_CLEANUP { AUTO | ON | OFF }
PROCESS_MAIN [ boolean ]
PROCESS_TOAST [ boolean ]
TRUNCATE [ boolean ]
PARALLEL integer
SKIP_DATABASE_STATS [ boolean ]
ONLY_DATABASE_STATS [ boolean ]
BUFFER_USAGE_LIMIT size
```

and *table_and_columns* is:

```
table_name [ ( column_name [, ...] ) ]
```

Description

VACUUM reclaims storage occupied by dead tuples. In normal Postgres Pro operation, tuples that are deleted or obsoleted by an update are not physically removed from their table; they remain present until a VACUUM is done. Therefore it's necessary to do VACUUM periodically, especially on frequently-updated tables.

Without a *table_and_columns* list, VACUUM processes every table and materialized view in the current database that the current user has permission to vacuum. With a list, VACUUM processes only those table(s).

VACUUM ANALYZE performs a VACUUM and then an ANALYZE for each selected table. This is a handy combination form for routine maintenance scripts. See [ANALYZE](#) for more details about its processing.

Plain VACUUM (without FULL) simply reclaims space and makes it available for re-use. This form of the command can operate in parallel with normal reading and writing of the table, as an exclusive lock is not obtained. However, extra space is not returned to the operating system (in most cases); it's just kept available for re-use within the same table. It also allows us to leverage multiple CPUs in order to process indexes. This feature is known as *parallel vacuum*. To disable this feature, one can use PARALLEL option and specify parallel workers as zero. VACUUM FULL rewrites the entire contents of the table into a new disk file with no extra space, allowing unused space to be returned to the operating system. This form is much slower and requires an ACCESS EXCLUSIVE lock on each table while it is being processed.

Parameters

FULL

Selects “full” vacuum, which can reclaim more space, but takes much longer and exclusively locks the table. This method also requires extra disk space, since it writes a new copy of the table and

doesn't release the old copy until the operation is complete. Usually this should only be used when a significant amount of space needs to be reclaimed from within the table.

FREEZE

Selects aggressive “freezing” of tuples. Specifying `FREEZE` is equivalent to performing `VACUUM` with the `vacuum_freeze_min_age` and `vacuum_freeze_table_age` parameters set to zero. Aggressive freezing is always performed when the table is rewritten, so this option is redundant when `FULL` is specified.

VERBOSE

Prints a detailed vacuum activity report for each table.

ANALYZE

Updates statistics used by the planner to determine the most efficient way to execute a query.

DISABLE_PAGE_SKIPPING

Normally, `VACUUM` will skip pages based on the [visibility map](#). Pages where all tuples are known to be frozen can always be skipped, and those where all tuples are known to be visible to all transactions may be skipped except when performing an aggressive vacuum. Furthermore, except when performing an aggressive vacuum, some pages may be skipped in order to avoid waiting for other sessions to finish using them. This option disables all page-skipping behavior, and is intended to be used only when the contents of the visibility map are suspect, which should happen only if there is a hardware or software issue causing database corruption.

SKIP_LOCKED

Specifies that `VACUUM` should not wait for any conflicting locks to be released when beginning work on a relation: if a relation cannot be locked immediately without waiting, the relation is skipped. Note that even with this option, `VACUUM` may still block when opening the relation's indexes. Additionally, `VACUUM ANALYZE` may still block when acquiring sample rows from partitions, table inheritance children, and some types of foreign tables. Also, while `VACUUM` ordinarily processes all partitions of specified partitioned tables, this option will cause `VACUUM` to skip all partitions if there is a conflicting lock on the partitioned table.

INDEX_CLEANUP

Normally, `VACUUM` will skip index vacuuming when there are very few dead tuples in the table. The cost of processing all of the table's indexes is expected to greatly exceed the benefit of removing dead index tuples when this happens. This option can be used to force `VACUUM` to process indexes when there are more than zero dead tuples. The default is `AUTO`, which allows `VACUUM` to skip index vacuuming when appropriate. If `INDEX_CLEANUP` is set to `ON`, `VACUUM` will conservatively remove all dead tuples from indexes. This may be useful for backwards compatibility with earlier releases of PostgreSQL where this was the standard behavior.

`INDEX_CLEANUP` can also be set to `OFF` to force `VACUUM` to *always* skip index vacuuming, even when there are many dead tuples in the table. This may be useful when it is necessary to make `VACUUM` run as quickly as possible to avoid imminent transaction ID wraparound (see [Section 24.1.5](#)). However, the wraparound failsafe mechanism controlled by `vacuum_failsafe_age` will generally trigger automatically to avoid transaction ID wraparound failure, and should be preferred. If index cleanup is not performed regularly, performance may suffer, because as the table is modified indexes will accumulate dead tuples and the table itself will accumulate dead line pointers that cannot be removed until index cleanup is completed.

This option has no effect for tables that have no index and is ignored if the `FULL` option is used. It also has no effect on the transaction ID wraparound failsafe mechanism. When triggered it will skip index vacuuming, even when `INDEX_CLEANUP` is set to `ON`.

PROCESS_MAIN

Specifies that `VACUUM` should attempt to process the main relation. This is usually the desired behavior and is the default. Setting this option to false may be useful when it is only necessary to vacuum a relation's corresponding `TOAST` table.

PROCESS_TOAST

Specifies that `VACUUM` should attempt to process the corresponding `TOAST` table for each relation, if one exists. This is usually the desired behavior and is the default. Setting this option to false may be useful when it is only necessary to vacuum the main relation. This option is required when the `FULL` option is used.

TRUNCATE

Specifies that `VACUUM` should attempt to truncate off any empty pages at the end of the table and allow the disk space for the truncated pages to be returned to the operating system. This is normally the desired behavior and is the default unless the `vacuum_truncate` option has been set to false for the table to be vacuumed. Setting this option to false may be useful to avoid `ACCESS EXCLUSIVE` lock on the table that the truncation requires. This option is ignored if the `FULL` option is used.

PARALLEL

Perform index vacuum and index cleanup phases of `VACUUM` in parallel using *integer* background workers (for the details of each vacuum phase, please refer to [Table 27.53](#)). The number of workers used to perform the operation is equal to the number of indexes on the relation that support parallel vacuum which is limited by the number of workers specified with `PARALLEL` option if any which is further limited by [max_parallel_maintenance_workers](#). An index can participate in parallel vacuum if and only if the size of the index is more than [min_parallel_index_scan_size](#). Please note that it is not guaranteed that the number of parallel workers specified in *integer* will be used during execution. It is possible for a vacuum to run with fewer workers than specified, or even with no workers at all. Only one worker can be used per index. So parallel workers are launched only when there are at least 2 indexes in the table. Workers for vacuum are launched before the start of each phase and exit at the end of the phase. These behaviors might change in a future release. This option can't be used with the `FULL` option.

SKIP_DATABASE_STATS

Specifies that `VACUUM` should skip updating the database-wide statistics about oldest unfrozen XIDs. Normally `VACUUM` will update these statistics once at the end of the command. However, this can take awhile in a database with a very large number of tables, and it will accomplish nothing unless the table that had contained the oldest unfrozen XID was among those vacuumed. Moreover, if multiple `VACUUM` commands are issued in parallel, only one of them can update the database-wide statistics at a time. Therefore, if an application intends to issue a series of many `VACUUM` commands, it can be helpful to set this option in all but the last such command; or set it in all the commands and separately issue `VACUUM (ONLY_DATABASE_STATS)` afterwards.

ONLY_DATABASE_STATS

Specifies that `VACUUM` should do nothing except update the database-wide statistics about oldest unfrozen XIDs. When this option is specified, the *table_and_columns* list must be empty, and no other option may be enabled except `VERBOSE`.

BUFFER_USAGE_LIMIT

Specifies the [Buffer Access Strategy](#) ring buffer size for `VACUUM`. This size is used to calculate the number of shared buffers which will be reused as part of this strategy. 0 disables use of a Buffer Access Strategy. If `ANALYZE` is also specified, the `BUFFER_USAGE_LIMIT` value is used for both the vacuum and analyze stages. This option can't be used with the `FULL` option except if `ANALYZE` is also specified. When this option is not specified, `VACUUM` uses the value from [vacuum_buffer_usage_limit](#). Higher settings can allow `VACUUM` to run more quickly, but having too large a setting may cause too

many other useful pages to be evicted from shared buffers. The minimum value is 128 kB and the maximum value is 16 GB.

boolean

Specifies whether the selected option should be turned on or off. You can write `TRUE`, `ON`, or `1` to enable the option, and `FALSE`, `OFF`, or `0` to disable it. The *boolean* value can also be omitted, in which case `TRUE` is assumed.

integer

Specifies a non-negative integer value passed to the selected option.

size

Specifies an amount of memory in kilobytes. Sizes may also be specified as a string containing the numerical size followed by any one of the following memory units: `B` (bytes), `kB` (kilobytes), `MB` (megabytes), `GB` (gigabytes), or `TB` (terabytes).

table_name

The name (optionally schema-qualified) of a specific table or materialized view to vacuum. If the specified table is a partitioned table, all of its leaf partitions are vacuumed.

column_name

The name of a specific column to analyze. Defaults to all columns. If a column list is specified, `ANALYZE` must also be specified.

Outputs

When `VERBOSE` is specified, `VACUUM` emits progress messages to indicate which table is currently being processed. Various statistics about the tables are printed as well.

Notes

To vacuum a table, one must ordinarily have the `MAINTAIN` privilege on the table. However, database owners are allowed to vacuum all tables in their databases, except shared catalogs. `VACUUM` will skip over any tables that the calling user does not have permission to vacuum.

While `VACUUM` is running, the [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`.

`VACUUM` cannot be executed inside a transaction block.

For tables with GIN indexes, `VACUUM` (in any form) also completes any pending index insertions, by moving pending index entries to the appropriate places in the main GIN index structure. See [Section 65.4.4.1](#) for details.

We recommend that all databases be vacuumed regularly in order to remove dead rows. Postgres Pro includes an “autovacuum” facility which can automate routine vacuum maintenance. For more information about automatic and manual vacuuming, see [Section 24.1](#).

The `FULL` option is not recommended for routine use, but might be useful in special cases. An example is when you have deleted or updated most of the rows in a table and would like the table to physically shrink to occupy less disk space and allow faster table scans. `VACUUM FULL` will usually shrink the table more than a plain `VACUUM` would.

The `PARALLEL` option is used only for vacuum purposes. If this option is specified with the `ANALYZE` option, it does not affect `ANALYZE`.

`VACUUM` causes a substantial increase in I/O traffic, which might cause poor performance for other active sessions. Therefore, it is sometimes advisable to use the cost-based vacuum delay feature. For parallel vacuum, each worker sleeps in proportion to the work done by that worker. See [Section 19.4.4](#) for details.

Each backend running `VACUUM` without the `FULL` option will report its progress in the `pg_stat_progress_vacuum` view. Backends running `VACUUM FULL` will instead report their progress in the `pg_stat_progress_cluster` view. See [Section 27.4.5](#) and [Section 27.4.2](#) for details.

Examples

To clean a single table `onek`, analyze it for the optimizer and print a detailed vacuum activity report:

```
VACUUM (VERBOSE, ANALYZE) onek;
```

Compatibility

There is no `VACUUM` statement in the SQL standard.

The following syntax was used before PostgreSQL version 9.0 and is still supported:

```
VACUUM [ FULL ] [ FREEZE ] [ VERBOSE ] [ ANALYZE ] [ table_and_columns [, ...] ]
```

Note that in this syntax, the options must be specified in exactly the order shown.

See Also

[vacuumdb](#), [Section 19.4.4](#), [Section 24.1.6](#), [Section 27.4.5](#), [Section 27.4.2](#)

VALUES

VALUES — compute a set of rows

Synopsis

```
VALUES ( expression [, ...] ) [, ...]  
  [ ORDER BY sort_expression [ ASC | DESC | USING operator ] [, ...] ]  
  [ LIMIT { count | ALL } ]  
  [ OFFSET start [ ROW | ROWS ] ]  
  [ FETCH { FIRST | NEXT } [ count ] { ROW | ROWS } ONLY ]
```

Description

VALUES computes a row value or set of row values specified by value expressions. It is most commonly used to generate a “constant table” within a larger command, but it can be used on its own.

When more than one row is specified, all the rows must have the same number of elements. The data types of the resulting table's columns are determined by combining the explicit or inferred types of the expressions appearing in that column, using the same rules as for UNION (see [Section 10.5](#)).

Within larger commands, VALUES is syntactically allowed anywhere that SELECT is. Because it is treated like a SELECT by the grammar, it is possible to use the ORDER BY, LIMIT (or equivalently FETCH FIRST), and OFFSET clauses with a VALUES command.

Parameters

expression

A constant or expression to compute and insert at the indicated place in the resulting table (set of rows). In a VALUES list appearing at the top level of an INSERT, an *expression* can be replaced by DEFAULT to indicate that the destination column's default value should be inserted. DEFAULT cannot be used when VALUES appears in other contexts.

sort_expression

An expression or integer constant indicating how to sort the result rows. This expression can refer to the columns of the VALUES result as column1, column2, etc. For more details see [ORDER BY Clause](#) in the [SELECT](#) documentation.

operator

A sorting operator. For details see [ORDER BY Clause](#) in the [SELECT](#) documentation.

count

The maximum number of rows to return. For details see [LIMIT Clause](#) in the [SELECT](#) documentation.

start

The number of rows to skip before starting to return rows. For details see [LIMIT Clause](#) in the [SELECT](#) documentation.

Notes

VALUES lists with very large numbers of rows should be avoided, as you might encounter out-of-memory failures or poor performance. VALUES appearing within INSERT is a special case (because the desired column types are known from the INSERT's target table, and need not be inferred by scanning the VALUES list), so it can handle larger lists than are practical in other contexts.

Examples

A bare VALUES command:

```
VALUES (1, 'one'), (2, 'two'), (3, 'three');
```

This will return a table of two columns and three rows. It's effectively equivalent to:

```
SELECT 1 AS column1, 'one' AS column2
UNION ALL
SELECT 2, 'two'
UNION ALL
SELECT 3, 'three';
```

More usually, VALUES is used within a larger SQL command. The most common use is in INSERT:

```
INSERT INTO films (code, title, did, date_prod, kind)
VALUES ('T_601', 'Yojimbo', 106, '1961-06-16', 'Drama');
```

In the context of INSERT, entries of a VALUES list can be DEFAULT to indicate that the column default should be used here instead of specifying a value:

```
INSERT INTO films VALUES
('UA502', 'Bananas', 105, DEFAULT, 'Comedy', '82 minutes'),
('T_601', 'Yojimbo', 106, DEFAULT, 'Drama', DEFAULT);
```

VALUES can also be used where a sub-SELECT might be written, for example in a FROM clause:

```
SELECT f.*
FROM films f, (VALUES('MGM', 'Horror'), ('UA', 'Sci-Fi')) AS t (studio, kind)
WHERE f.studio = t.studio AND f.kind = t.kind;
```

```
UPDATE employees SET salary = salary * v.increase
FROM (VALUES(1, 200000, 1.2), (2, 400000, 1.4)) AS v (depno, target, increase)
WHERE employees.depno = v.depno AND employees.sales >= v.target;
```

Note that an AS clause is required when VALUES is used in a FROM clause, just as is true for SELECT. It is not required that the AS clause specify names for all the columns, but it's good practice to do so. (The default column names for VALUES are column1, column2, etc. in Postgres Pro, but these names might be different in other database systems.)

When VALUES is used in INSERT, the values are all automatically coerced to the data type of the corresponding destination column. When it's used in other contexts, it might be necessary to specify the correct data type. If the entries are all quoted literal constants, coercing the first is sufficient to determine the assumed type for all:

```
SELECT * FROM machines
WHERE ip_address IN (VALUES('192.168.0.1'::inet), ('192.168.0.10'), ('192.168.1.43'));
```

Tip

For simple IN tests, it's better to rely on the [list-of-scalars](#) form of IN than to write a VALUES query as shown above. The list of scalars method requires less writing and is often more efficient.

Compatibility

VALUES conforms to the SQL standard. LIMIT and OFFSET are Postgres Pro extensions; see also under [SELECT](#).

See Also

[INSERT](#), [SELECT](#)

Postgres Pro Client Applications

This part contains reference information for Postgres Pro client applications and utilities. Not all of these commands are of general utility; some might require special privileges. The common feature of these applications is that they can be run on any host, independent of where the database server resides.

When specified on the command line, user and database names have their case preserved — the presence of spaces or special characters might require quoting. Table names and other identifiers do not have their case preserved, except where documented, and might require quoting.

clusterdb

clusterdb — cluster a Postgres Pro database

Synopsis

```
clusterdb [connection-option...] [option...] [ --table | -t table ] ... [ dbname | -a | --all ]
```

Description

clusterdb is a utility for recluster tables in a Postgres Pro database. It finds tables that have previously been clustered, and clusters them again on the same index that was last used. Tables that have never been clustered are not affected.

clusterdb is a wrapper around the SQL command [CLUSTER](#). There is no effective difference between clustering databases via this utility and via other methods for accessing the server.

Options

clusterdb accepts the following command-line arguments:

-a
--all

Cluster all databases.

[-d] dbname
[--dbname=] dbname

Specifies the name of the database to be clustered, when -a/--all is not used. If this is not specified, the database name is read from the environment variable PGDATABASE. If that is not set, the user name specified for the connection is used. The dbname can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

-e
--echo

Echo the commands that clusterdb generates and sends to the server.

-q
--quiet

Do not display progress messages.

-t table
--table=table

Cluster table only. Multiple tables can be clustered by writing multiple -t switches.

-v
--verbose

Print detailed information during processing.

-V
--version

Print the clusterdb version and exit.

-?
--help

Show help about clusterdb command line arguments, and exit.

clusterdb also accepts the following command-line arguments for connection parameters:

`-h host`
`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`
`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U username`
`--username=username`

User name to connect as.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force clusterdb to prompt for a password before connecting to a database.

This option is never essential, since clusterdb will automatically prompt for a password if the server demands password authentication. However, clusterdb will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`--maintenance-db=dbname`

When the `-a/--all` is used, connect to this database to gather the list of databases to cluster. If not specified, the `postgres` database will be used, or if that does not exist, `template1` will be used. This can be a [connection string](#). If so, connection string parameters will override any conflicting command line options. Also, connection string parameters other than the database name itself will be re-used when connecting to other databases.

Environment

PGDATABASE
PGHOST
PGPORT
PGUSER

Default connection parameters

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [CLUSTER](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To cluster the database `test`:

```
$ clusterdb test
```

To cluster a single table `foo` in a database named `xyzzzy`:

```
$ clusterdb --table=foo xyzzzy
```

See Also

[CLUSTER](#)

createdb

createdb — create a new Postgres Pro database

Synopsis

```
createdb [connection-option...] [option...] [dbname [description]]
```

Description

createdb creates a new Postgres Pro database.

Note

Note createdb only works for local DBs.

Normally, the database user who executes this command becomes the owner of the new database. However, a different owner can be specified via the `-O` option, if the executing user has appropriate privileges.

createdb is a wrapper around the SQL command `CREATE DATABASE`. There is no effective difference between creating databases via this utility and via other methods for accessing the server.

Options

createdb accepts the following command-line arguments:

dbname

Specifies the name of the database to be created. The name must be unique among all Postgres Pro databases in this cluster. The default is to create a database with the same name as the current system user.

description

Specifies a comment to be associated with the newly created database.

`-D tablespace`

`--tablespace=tablespace`

Specifies the default tablespace for the database. (This name is processed as a double-quoted identifier.)

`-e`

`--echo`

Echo the commands that createdb generates and sends to the server.

`-E encoding`

`--encoding=encoding`

Specifies the character encoding scheme to be used in this database. The character sets supported by the Postgres Pro server are described in [Section 23.3.1](#).

`-l locale`

`--locale=locale`

Specifies the locale to be used in this database. This is equivalent to specifying `--lc-collate`, `--lc-ctype`, and `--icu-locale` to the same value. Some locales are only valid for ICU and must be set with `--icu-locale`.

`--lc-collate=locale`

Specifies the LC_COLLATE setting to be used in this database.

`--lc-ctype=locale`

Specifies the LC_CTYPE setting to be used in this database.

`--builtin-locale=locale`

Specifies the locale name when the builtin provider is used. Locale support is described in [Section 23.1](#).

`--icu-locale=locale`

Specifies the ICU locale ID to be used in this database, if the ICU locale provider is selected.

`--icu-rules=rules`

Specifies additional collation rules to customize the behavior of the default collation of this database. This is supported for ICU only.

`--locale-provider={builtin|libc|icu}`

Specifies the locale provider for the database's default collation.

`-O owner`

`--owner=owner`

Specifies the database user who will own the new database. (This name is processed as a double-quoted identifier.)

`-S strategy`

`--strategy=strategy`

Specifies the database creation strategy. See [CREATE DATABASE STRATEGY](#) for more details.

`-T template`

`--template=template`

Specifies the template database from which to build this database. (This name is processed as a double-quoted identifier.)

`-V`

`--version`

Print the createdb version and exit.

`-?`

`--help`

Show help about createdb command line arguments, and exit.

The options `-D`, `-l`, `-E`, `-O`, and `-T` correspond to options of the underlying SQL command [CREATE DATABASE](#); see there for more information about them.

createdb also accepts the following command-line arguments for connection parameters:

`-h host`

`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`

`--port=port`

Specifies the TCP port or the local Unix domain socket file extension on which the server is listening for connections.

```
-U username
--username=username
```

User name to connect as.

```
-w
--no-password
```

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

```
-W
--password
```

Force createdb to prompt for a password before connecting to a database.

This option is never essential, since createdb will automatically prompt for a password if the server demands password authentication. However, createdb will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

```
--maintenance-db=dbname
```

Specifies the name of the database to connect to when creating the new database. If not specified, the `postgres` database will be used; if that does not exist (or if it is the name of the new database being created), `template1` will be used. This can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

Environment

`PGDATABASE`

If set, the name of the database to create, unless overridden on the command line.

`PGHOST`
`PGPORT`
`PGUSER`

Default connection parameters. `PGUSER` also determines the name of the database to create, if it is not specified on the command line or by `PGDATABASE`.

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [CREATE DATABASE](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To create the database `demo` using the default database server:

```
$ createdb demo
```

To create the database `demo` using the server on host `eden`, port 5000, using the `template0` template database, here is the command-line command and the underlying SQL command:

```
$ createdb -p 5000 -h eden -T template0 -e demo  
CREATE DATABASE demo TEMPLATE template0;
```

See Also

[dropdb](#), [CREATE DATABASE](#)

createuser

createuser — define a new Postgres Pro user account

Synopsis

```
createuser [connection-option...] [option...] [username]
```

Description

createuser creates a new Postgres Pro user (or more precisely, a role). Only superusers and users with `CREATEROLE` privilege can create new users, so createuser must be invoked by someone who can connect as a superuser or a user with `CREATEROLE` privilege.

Note

Note createuser cannot create global users.

If you wish to create a role with the `SUPERUSER`, `REPLICATION`, or `BYPASSRLS` privilege, you must connect as a superuser, not merely with `CREATEROLE` privilege. Being a superuser implies the ability to bypass all access permission checks within the database, so superuser access should not be granted lightly. `CREATEROLE` also conveys [very extensive privileges](#).

createuser is a wrapper around the SQL command `CREATE ROLE`. There is no effective difference between creating users via this utility and via other methods for accessing the server.

Options

createuser accepts the following command-line arguments:

username

Specifies the name of the Postgres Pro user to be created. This name must be different from all existing roles in this Postgres Pro installation.

`-a role`

`--with-admin=role`

Specifies an existing role that will be automatically added as a member of the new role with admin option, giving it the right to grant membership in the new role to others. Multiple existing roles can be specified by writing multiple `-a` switches.

`-c number`

`--connection-limit=number`

Set a maximum number of connections for the new user. The default is to set no limit.

`-d`

`--createdb`

The new user will be allowed to create databases.

`-D`

`--no-createdb`

The new user will not be allowed to create databases. This is the default.

`-e`

`--echo`

Echo the commands that createuser generates and sends to the server.

`-E`
`--encrypted`

This option is obsolete but still accepted for backward compatibility.

`-g role`
`--member-of=role`
`--role=role` (deprecated)

Specifies the new role should be automatically added as a member of the specified existing role. Multiple existing roles can be specified by writing multiple `-g` switches.

`-i`
`--inherit`

The new role will automatically inherit privileges of roles it is a member of. This is the default.

`-I`
`--no-inherit`

The new role will not automatically inherit privileges of roles it is a member of.

`--interactive`

Prompt for the user name if none is specified on the command line, and also prompt for whichever of the options `-d/-D`, `-r/-R`, `-s/-S` is not specified on the command line. (This was the default behavior up to PostgreSQL 9.1.)

`-l`
`--login`

The new user will be allowed to log in (that is, the user name can be used as the initial session user identifier). This is the default.

`-L`
`--no-login`

The new user will not be allowed to log in. (A role without login privilege is still useful as a means of managing database permissions.)

`-m role`
`--with-member=role`

Specifies an existing role that will be automatically added as a member of the new role. Multiple existing roles can be specified by writing multiple `-m` switches.

`-P`
`--pwprompt`

If given, `createuser` will issue a prompt for the password of the new user. This is not necessary if you do not plan on using password authentication.

`-r`
`--createrole`

The new user will be allowed to create, alter, drop, comment on, change the security label for other roles; that is, this user will have `CREATEROLE` privilege. See [role creation](#) for more details about what capabilities are conferred by this privilege.

`-R`
`--no-createrole`

The new user will not be allowed to create new roles. This is the default.

`-s`
`--superuser`
The new user will be a superuser.

`-S`
`--no-superuser`
The new user will not be a superuser. This is the default.

`-v timestamp`
`--valid-until=timestamp`
Set a date and time after which the role's password is no longer valid. The default is to set no password expiry date.

`-V`
`--version`
Print the createuser version and exit.

`--bypassrls`
The new user will bypass every row-level security (RLS) policy.

`--no-bypassrls`
The new user will not bypass row-level security (RLS) policies. This is the default.

`--profile=profile`
The new role will be assigned the profile that enforces password management policy. All restrictions of this profile will apply to the role.

`--replication`
The new user will have the `REPLICATION` privilege, which is described more fully in the documentation for [CREATE ROLE](#).

`--no-replication`
The new user will not have the `REPLICATION` privilege, which is described more fully in the documentation for [CREATE ROLE](#). This is the default.

`-?`
`--help`
Show help about createuser command line arguments, and exit.

createuser also accepts the following command-line arguments for connection parameters:

`-h host`
`--host=host`
Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`
`--port=port`
Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U username`
`--username=username`
User name to connect as (not the user name to create).

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force createuser to prompt for a password (for connecting to the server, not for the password of the new user).

This option is never essential, since createuser will automatically prompt for a password if the server demands password authentication. However, createuser will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

Environment

`PGHOST`
`PGPORT`
`PGUSER`

Default connection parameters

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [CREATE ROLE](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To create a user `joe` on the default database server:

```
$ createuser joe
```

To create a user `joe` on the default database server with prompting for some additional attributes:

```
$ createuser --interactive joe
Shall the new role be a superuser? (y/n) n
Shall the new role be allowed to create databases? (y/n) n
Shall the new role be allowed to create more new roles? (y/n) n
```

To create the same user `joe` using the server on host `eden`, port 5000, with attributes explicitly specified, taking a look at the underlying command:

```
$ createuser -h eden -p 5000 -S -D -R -e joe
CREATE ROLE joe NOSUPERUSER NOCREATEDB NOCREATEROLE INHERIT LOGIN;
```

To create the user `joe` as a superuser, and assign a password immediately:

```
$ createuser -P -s -e joe
Enter password for new role: xyzzy
Enter it again: xyzzy
```

```
CREATE ROLE joe PASSWORD 'md5b5f5ba1a423792b526f799ae4eb3d59e' SUPERUSER CREATEDB  
CREATEROLE INHERIT LOGIN;
```

In the above example, the new password isn't actually echoed when typed, but we show what was typed for clarity. As you see, the password is encrypted before it is sent to the client.

See Also

[dropuser](#), [CREATE ROLE](#), [createrole_self_grant](#)

dropdb

dropdb — remove a Postgres Pro database

Synopsis

```
dropdb [connection-option...] [option...] dbname
```

Description

dropdb destroys an existing Postgres Pro database. The user who executes this command must be a database superuser or the owner of the database.

dropdb is a wrapper around the SQL command [DROP DATABASE](#). There is no effective difference between dropping databases via this utility and via other methods for accessing the server.

Options

dropdb accepts the following command-line arguments:

dbname

Specifies the name of the database to be removed.

-e

--echo

Echo the commands that dropdb generates and sends to the server.

-f

--force

Attempt to terminate all existing connections to the target database before dropping it. See [DROP DATABASE](#) for more information on this option.

-i

--interactive

Issues a verification prompt before doing anything destructive.

-V

--version

Print the dropdb version and exit.

--if-exists

Do not throw an error if the database does not exist. A notice is issued in this case.

-?

--help

Show help about dropdb command line arguments, and exit.

dropdb also accepts the following command-line arguments for connection parameters:

-h *host*

--host=*host*

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`
`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U username`
`--username=username`

User name to connect as.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force dropdb to prompt for a password before connecting to a database.

This option is never essential, since dropdb will automatically prompt for a password if the server demands password authentication. However, dropdb will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`--maintenance-db=dbname`

Specifies the name of the database to connect to in order to drop the target database. If not specified, the `postgres` database will be used; if that does not exist (or is the database being dropped), `template1` will be used. This can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

Environment

PGHOST
PGPORT
PGUSER

Default connection parameters

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [DROP DATABASE](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To destroy the database `demo` on the default database server:

```
$ dropdb demo
```

To destroy the database `demo` using the server on host `eden`, port 5000, with verification and a peek at the underlying command:

```
$ dropdb -p 5000 -h eden -i -e demo
Database "demo" will be permanently deleted.
Are you sure? (y/n) y
DROP DATABASE demo;
```

See Also

[createdb](#), [DROP DATABASE](#)

dropuser

dropuser — remove a Postgres Pro user account

Synopsis

```
dropuser [connection-option...] [option...] [username]
```

Description

dropuser removes an existing Postgres Pro user. Superusers can use this command to remove any role; otherwise, only non-superuser roles can be removed, and only by a user who possesses the `CREATEROLE` privilege and has been granted `ADMIN OPTION` on the target role.

dropuser is a wrapper around the SQL command `DROP ROLE`. There is no effective difference between dropping users via this utility and via other methods for accessing the server.

Options

dropuser accepts the following command-line arguments:

username

Specifies the name of the Postgres Pro user to be removed. You will be prompted for a name if none is specified on the command line and the `-i/--interactive` option is used.

`-e`

`--echo`

Echo the commands that dropuser generates and sends to the server.

`-i`

`--interactive`

Prompt for confirmation before actually removing the user, and prompt for the user name if none is specified on the command line.

`-V`

`--version`

Print the dropuser version and exit.

`--if-exists`

Do not throw an error if the user does not exist. A notice is issued in this case.

`-?`

`--help`

Show help about dropuser command line arguments, and exit.

dropuser also accepts the following command-line arguments for connection parameters:

`-h host`

`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`

`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

```
-U username
--username=username
```

User name to connect as (not the user name to drop).

```
-w
--no-password
```

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

```
-W
--password
```

Force dropuser to prompt for a password before connecting to a database.

This option is never essential, since dropuser will automatically prompt for a password if the server demands password authentication. However, dropuser will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

Environment

```
PGHOST
PGPORT
PGUSER
```

Default connection parameters

```
PG_COLOR
```

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [DROP ROLE](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To remove user `joe` from the default database server:

```
$ dropuser joe
```

To remove user `joe` using the server on host `eden`, port 5000, with verification and a peek at the underlying command:

```
$ dropuser -p 5000 -h eden -i -e joe
Role "joe" will be permanently removed.
Are you sure? (y/n) y
DROP ROLE joe;
```

See Also

[createuser](#), [DROP ROLE](#)

ecpg

ecpg — embedded SQL C preprocessor

Synopsis

`ecpg [option...] file...`

Description

`ecpg` is the embedded SQL preprocessor for C programs. It converts C programs with embedded SQL statements to normal C code by replacing the SQL invocations with special function calls. The output files can then be processed with any C compiler tool chain.

`ecpg` will convert each input file given on the command line to the corresponding C output file. If an input file name does not have any extension, `.pgc` is assumed. The file's extension will be replaced by `.c` to construct the output file name. But the output file name can be overridden using the `-o` option.

If an input file name is just `-`, `ecpg` reads the program from standard input (and writes to standard output, unless that is overridden with `-o`).

This reference page does not describe the embedded SQL language. See [Chapter 37](#) for more information on that topic.

Options

`ecpg` accepts the following command-line arguments:

- `-c`
Automatically generate certain C code from SQL code. Currently, this works for `EXEC SQL TYPE`.
- `-C mode`
Set a compatibility mode. *mode* can be `INFORMIX`, `INFORMIX_SE`, or `ORACLE`.
- `-D symbol[=value]`
Define a preprocessor symbol, equivalently to the `EXEC SQL DEFINE` directive. If no *value* is specified, the symbol is defined with the value `1`.
- `-h`
Process header files. When this option is specified, the output file extension becomes `.h` not `.c`, and the default input file extension is `.pgh` not `.pgc`. Also, the `-c` option is forced on.
- `-i`
Parse system include files as well.
- `-I directory`
Specify an additional include path, used to find files included via `EXEC SQL INCLUDE`. Defaults are `.` (current directory), `/usr/local/include`, the Postgres Pro include directory which is defined at compile time (default: `/usr/local/pgsql/include`), and `/usr/include`, in that order.
- `-o filename`
Specifies that `ecpg` should write all its output to the given *filename*. Write `-o -` to send all output to standard output.
- `-r option`
Selects run-time behavior. *Option* can be one of the following:

no_indicator

Do not use indicators but instead use special values to represent null values. Historically there have been databases using this approach.

prepare

Prepare all statements before using them. Libecpg will keep a cache of prepared statements and reuse a statement if it gets executed again. If the cache runs full, libecpg will free the least used statement.

questionmarks

Allow question mark as placeholder for compatibility reasons. This used to be the default long ago.

-t

Turn on autocommit of transactions. In this mode, each SQL command is automatically committed unless it is inside an explicit transaction block. In the default mode, commands are committed only when EXEC SQL COMMIT is issued.

-v

Print additional information including the version and the "include" path.

--version

Print the ecpg version and exit.

-?**--help**

Show help about ecpg command line arguments, and exit.

Notes

When compiling the preprocessed C code files, the compiler needs to be able to find the ECPG header files in the Postgres Pro include directory. Therefore, you might have to use the `-I` option when invoking the compiler (e.g., `-I/usr/local/pgsql/include`).

Programs using C code with embedded SQL have to be linked against the `libecpg` library, for example using the linker options `-L/usr/local/pgsql/lib -lecpg`.

The value of either of these directories that is appropriate for the installation can be found out using [pg_config](#).

Examples

If you have an embedded SQL C source file named `prog1.pgc`, you can create an executable program using the following sequence of commands:

```
ecpg prog1.pgc
cc -I/usr/local/pgsql/include -c prog1.c
cc -o prog1 prog1.o -L/usr/local/pgsql/lib -lecpg
```

pg_amcheck

pg_amcheck — checks for corruption in one or more Postgres Pro databases

Synopsis

```
pg_amcheck [option...] [dbname]
```

Description

pg_amcheck supports running [amcheck](#)'s corruption checking functions against one or more databases, with options to select which schemas, tables and indexes to check, which kinds of checking to perform, and whether to perform the checks in parallel, and if so, the number of parallel connections to establish and use.

Only ordinary and toast table relations, materialized views, sequences, and btree indexes are currently supported. Other relation types are silently skipped.

If *dbname* is specified, it should be the name of a single database to check, and no other database selection options should be present. Otherwise, if any database selection options are present, all matching databases will be checked. If no such options are present, the default database will be checked. Database selection options include `--all`, `--database` and `--exclude-database`. They also include `--relation`, `--exclude-relation`, `--table`, `--exclude-table`, `--index`, and `--exclude-index`, but only when such options are used with a three-part pattern (e.g. `mydb*.myschema*.myrel*`). Finally, they include `--schema` and `--exclude-schema` when such options are used with a two-part pattern (e.g. `mydb*.myschema*`).

dbname can also be a [connection string](#).

Options

The following command-line options control what is checked:

```
-a  
--all
```

Check all databases, except for any excluded via `--exclude-database`.

```
-d pattern  
--database=pattern
```

Check databases matching the specified *pattern*, except for any excluded by `--exclude-database`. This option can be specified more than once.

```
-D pattern  
--exclude-database=pattern
```

Exclude databases matching the given *pattern*. This option can be specified more than once.

```
-i pattern  
--index=pattern
```

Check indexes matching the specified *pattern*, unless they are otherwise excluded. This option can be specified more than once.

This is similar to the `--relation` option, except that it applies only to indexes, not to other relation types.

```
-I pattern  
--exclude-index=pattern
```

Exclude indexes matching the specified *pattern*. This option can be specified more than once.

This is similar to the `--exclude-relation` option, except that it applies only to indexes, not other relation types.

```
-r pattern  
--relation=pattern
```

Check relations matching the specified *pattern*, unless they are otherwise excluded. This option can be specified more than once.

Patterns may be unqualified, e.g. `myrel*`, or they may be schema-qualified, e.g. `myschema*.myrel*` or database-qualified and schema-qualified, e.g. `mydb*.myschema*.myrel*`. A database-qualified pattern will add matching databases to the list of databases to be checked.

```
-R pattern  
--exclude-relation=pattern
```

Exclude relations matching the specified *pattern*. This option can be specified more than once.

As with `--relation`, the *pattern* may be unqualified, schema-qualified, or database- and schema-qualified.

```
-s pattern  
--schema=pattern
```

Check tables and indexes in schemas matching the specified *pattern*, unless they are otherwise excluded. This option can be specified more than once.

To select only tables in schemas matching a particular pattern, consider using something like `--table=SCHEMAPAT.* --no-dependent-indexes`. To select only indexes, consider using something like `--index=SCHEMAPAT.*`.

A schema pattern may be database-qualified. For example, you may write `--schema=mydb*.myschema*` to select schemas matching `myschema*` in databases matching `mydb*`.

```
-S pattern  
--exclude-schema=pattern
```

Exclude tables and indexes in schemas matching the specified *pattern*. This option can be specified more than once.

As with `--schema`, the pattern may be database-qualified.

```
-t pattern  
--table=pattern
```

Check tables matching the specified *pattern*, unless they are otherwise excluded. This option can be specified more than once.

This is similar to the `--relation` option, except that it applies only to tables, materialized views, and sequences, not to indexes.

```
-T pattern  
--exclude-table=pattern
```

Exclude tables matching the specified *pattern*. This option can be specified more than once.

This is similar to the `--exclude-relation` option, except that it applies only to tables, materialized views, and sequences, not to indexes.

```
--no-dependent-indexes
```

By default, if a table is checked, any btree indexes of that table will also be checked, even if they are not explicitly selected by an option such as `--index` or `--relation`. This option suppresses that behavior.

`--no-dependent-toast`

By default, if a table is checked, its toast table, if any, will also be checked, even if it is not explicitly selected by an option such as `--table` or `--relation`. This option suppresses that behavior.

`--no-strict-names`

By default, if an argument to `--database`, `--table`, `--index`, or `--relation` matches no objects, it is a fatal error. This option downgrades that error to a warning.

The following command-line options control checking of tables:

`--exclude-toast-pointers`

By default, whenever a toast pointer is encountered in a table, a lookup is performed to ensure that it references apparently-valid entries in the toast table. These checks can be quite slow, and this option can be used to skip them.

`--on-error-stop`

After reporting all corruptions on the first page of a table where corruption is found, stop processing that table relation and move on to the next table or index.

Note that index checking always stops after the first corrupt page. This option only has meaning relative to table relations.

`--skip=option`

If `all-frozen` is given, table corruption checks will skip over pages in all tables that are marked as all frozen.

If `all-visible` is given, table corruption checks will skip over pages in all tables that are marked as all visible.

By default, no pages are skipped. This can be specified as `none`, but since this is the default, it need not be mentioned.

`--startblock=block`

Start checking at the specified block number. An error will occur if the table relation being checked has fewer than this number of blocks. This option does not apply to indexes, and is probably only useful when checking a single table relation. See `--endblock` for further caveats.

`--endblock=block`

End checking at the specified block number. An error will occur if the table relation being checked has fewer than this number of blocks. This option does not apply to indexes, and is probably only useful when checking a single table relation. If both a regular table and a toast table are checked, this option will apply to both, but higher-numbered toast blocks may still be accessed while validating toast pointers, unless that is suppressed using `--exclude-toast-pointers`.

The following command-line options control checking of B-tree indexes:

`--checkunique`

For each index with unique constraint checked, verify that no more than one among duplicate entries is visible in the index using [amcheck's](#) `checkunique` option.

`--heapallindexed`

For each index checked, verify the presence of all heap tuples as index tuples in the index using [amcheck's](#) `heapallindexed` option.

`--parent-check`

For each btree index checked, use [amcheck's](#) `bt_index_parent_check` function, which performs additional checks of parent/child relationships during index checking.

The default is to use `amcheck`'s `bt_index_check` function, but note that use of the `--rootdescend` option implicitly selects `bt_index_parent_check`.

`--rootdescend`

For each index checked, re-find tuples on the leaf level by performing a new search from the root page for each tuple using `amcheck`'s `rootdescend` option.

Use of this option implicitly also selects the `--parent-check` option.

This form of verification was originally written to help in the development of btree index features. It may be of limited use or even of no use in helping detect the kinds of corruption that occur in practice. It may also cause corruption checking to take considerably longer and consume considerably more resources on the server.

Warning

The extra checks performed against B-tree indexes when the `--parent-check` option or the `--rootdescend` option is specified require relatively strong relation-level locks. These checks are the only checks that will block concurrent data modification from `INSERT`, `UPDATE`, and `DELETE` commands.

The following command-line options control the connection to the server:

`-h hostname`

`--host=hostname`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`

`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U`

`--username=username`

User name to connect as.

`-w`

`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`

`--password`

Force `pg_amcheck` to prompt for a password before connecting to a database.

This option is never essential, since `pg_amcheck` will automatically prompt for a password if the server demands password authentication. However, `pg_amcheck` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`--maintenance-db=dbname`

Specifies a database or [connection string](#) to be used to discover the list of databases to be checked. If neither `--all` nor any option including a database pattern is used, no such connection is required and

this option does nothing. Otherwise, any connection string parameters other than the database name which are included in the value for this option will also be used when connecting to the databases being checked. If this option is omitted, the default is `postgres` or, if that fails, `template1`.

Other options are also available:

`-e`
`--echo`

Echo to stdout all SQL sent to the server.

`-j num`
`--jobs=num`

Use *num* concurrent connections to the server, or one per object to be checked, whichever is less.

The default is to use a single connection.

`-P`
`--progress`

Show progress information. Progress information includes the number of relations for which checking has been completed, and the total size of those relations. It also includes the total number of relations that will eventually be checked, and the estimated size of those relations.

`-v`
`--verbose`

Print more messages. In particular, this will print a message for each relation being checked, and will increase the level of detail shown for server errors.

`-V`
`--version`

Print the `pg_amcheck` version and exit.

`--install-missing`
`--install-missing=schema`

Install any missing extensions that are required to check the database(s). If not yet installed, each extension's objects will be installed into the given *schema*, or if not specified into schema `pg_catalog`.

At present, the only required extension is [amcheck](#).

`-?`
`--help`

Show help about `pg_amcheck` command line arguments, and exit.

Environment

`pg_amcheck`, like most other Postgres Pro utilities, also uses the environment variables supported by `libpq` (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

`pg_amcheck` is designed to work with Postgres Pro 14.0 and later.

See Also

[amcheck](#)

pgbench

pgbench — run a benchmark test on Postgres Pro

Synopsis

```
pgbench -i [option...] [dbname]
```

```
pgbench [option...] [dbname]
```

Description

pgbench is a simple program for running benchmark tests on Postgres Pro. It runs the same sequence of SQL commands over and over, possibly in multiple concurrent database sessions, and then calculates the average transaction rate (transactions per second). By default, pgbench tests a scenario that is loosely based on TPC-B, involving five `SELECT`, `UPDATE`, and `INSERT` commands per transaction. However, it is easy to test other cases by writing your own transaction script files.

Typical output from pgbench looks like:

```
transaction type: <builtin: TPC-B (sort of)>
scaling factor: 10
query mode: simple
number of clients: 10
number of threads: 1
maximum number of tries: 1
number of transactions per client: 1000
number of transactions actually processed: 10000/10000
number of failed transactions: 0 (0.000%)
latency average = 11.013 ms
latency stddev = 7.351 ms
initial connection time = 45.758 ms
tps = 896.967014 (without initial connection time)
```

The first seven lines report some of the most important parameter settings. The sixth line reports the maximum number of tries for transactions with serialization or deadlock errors (see [Failures and Serialization/Deadlock Retries](#) for more information). The eighth line reports the number of transactions completed and intended (the latter being just the product of number of clients and number of transactions per client); these will be equal unless the run failed before completion or some SQL command(s) failed. (In `-T` mode, only the actual number of transactions is printed.) The next line reports the number of failed transactions due to serialization or deadlock errors (see [Failures and Serialization/Deadlock Retries](#) for more information). The last line reports the number of transactions per second.

The default TPC-B-like transaction test requires specific tables to be set up beforehand. pgbench should be invoked with the `-i` (initialize) option to create and populate these tables. (When you are testing a custom script, you don't need this step, but will instead need to do whatever setup your test needs.) Initialization looks like:

```
pgbench -i [ other-options ] dbname
```

where *dbname* is the name of the already-created database to test in. (You may also need `-h`, `-p`, and/or `-U` options to specify how to connect to the database server.)

Caution

pgbench -i creates four tables `pgbench_accounts`, `pgbench_branches`, `pgbench_history`, and `pgbench_tellers`, destroying any existing tables of these names. Be very careful to use another database if you have tables having these names!

At the default “scale factor” of 1, the tables initially contain this many rows:

table	# of rows
-----	-----
pgbench_branches	1
pgbench_tellers	10
pgbench_accounts	100000
pgbench_history	0

You can (and, for most purposes, probably should) increase the number of rows by using the `-s` (scale factor) option. The `-F` (fillfactor) option might also be used at this point.

Once you have done the necessary setup, you can run your benchmark with a command that doesn't include `-i`, that is

```
pgbench [ options ] dbname
```

In nearly all cases, you'll need some options to make a useful test. The most important options are `-c` (number of clients), `-t` (number of transactions), `-T` (time limit), and `-f` (specify a custom script file). See below for a full list.

Options

The following is divided into three subsections. Different options are used during database initialization and while running benchmarks, but some options are useful in both cases.

Initialization Options

pgbench accepts the following command-line initialization arguments:

```
[-d] dbname  
[--dbname=] dbname
```

Specifies the name of the database to test in. If this is not specified, the environment variable `PGDATABASE` is used. If that is not set, the user name specified for the connection is used.

```
-i  
--initialize
```

Required to invoke initialization mode.

```
-I init_steps  
--init-steps=init_steps
```

Perform just a selected set of the normal initialization steps. *init_steps* specifies the initialization steps to be performed, using one character per step. Each step is invoked in the specified order. The default is `dtgvp`. The available steps are:

d (Drop)

Drop any existing pgbench tables.

t (create Tables)

Create the tables used by the standard pgbench scenario, namely `pgbench_accounts`, `pgbench_branches`, `pgbench_history`, and `pgbench_tellers`.

g or G (Generate data, client-side or server-side)

Generate data and load it into the standard tables, replacing any data already present.

With `g` (client-side data generation), data is generated in `pgbench` client and then sent to the server. This uses the client/server bandwidth extensively through a `COPY`. `pgbench` uses the `FREEZE`

option to load data into ordinary (non-partition) tables with version 14 or later of Postgres Pro to speed up subsequent `VACUUM`. Using `g` causes logging to print one message every 100,000 rows while generating data for all tables.

With `G` (server-side data generation), only small queries are sent from the `pgbench` client and then data is actually generated in the server. No significant bandwidth is required for this variant, but the server will do more work. Using `G` causes logging not to print any progress message while generating data.

The default initialization behavior uses client-side data generation (equivalent to `g`).

`v` (Vacuum)

Invoke `VACUUM` on the standard tables.

`p` (create Primary keys)

Create primary key indexes on the standard tables.

`f` (create Foreign keys)

Create foreign key constraints between the standard tables. (Note that this step is not performed by default.)

`-F fillfactor`

`--fillfactor=fillfactor`

Create the `pgbench_accounts`, `pgbench_tellers` and `pgbench_branches` tables with the given fillfactor. Default is 100.

`-n`

`--no-vacuum`

Perform no vacuuming during initialization. (This option suppresses the `v` initialization step, even if it was specified in `-I`.)

`-q`

`--quiet`

Switch logging to quiet mode, producing only one progress message per 5 seconds. The default logging prints one message each 100,000 rows, which often outputs many lines per second (especially on good hardware).

This setting has no effect if `G` is specified in `-I`.

`-s scale_factor`

`--scale=scale_factor`

Multiply the number of rows generated by the scale factor. For example, `-s 100` will create 10,000,000 rows in the `pgbench_accounts` table. Default is 1. When the scale is 20,000 or larger, the columns used to hold account identifiers (`aid` columns) will switch to using larger integers (`bigint`), in order to be big enough to hold the range of account identifiers.

`--foreign-keys`

Create foreign key constraints between the standard tables. (This option adds the `f` step to the initialization step sequence, if it is not already present.)

`--index-tablespace=index_tablespace`

Create indexes in the specified tablespace, rather than the default tablespace.

`--partition-method=NAME`

Create a partitioned `pgbench_accounts` table with *NAME* method. Expected values are `range` or `hash`. This option requires that `--partitions` is set to non-zero. If unspecified, default is `range`.

`--partitions=NUM`

Create a partitioned `pgbench_accounts` table with *NUM* partitions of nearly equal size for the scaled number of accounts. Default is 0, meaning no partitioning.

`--tablespace=tablespace`

Create tables in the specified tablespace, rather than the default tablespace.

`--unlogged-tables`

Create all tables as unlogged tables, rather than permanent tables.

Benchmarking Options

`pgbench` accepts the following command-line benchmarking arguments:

`-b scriptname[@weight]`

`--builtin=scriptname[@weight]`

Add the specified built-in script to the list of scripts to be executed. Available built-in scripts are: `tpcb-like`, `simple-update` and `select-only`. Unambiguous prefixes of built-in names are accepted. With the special name `list`, show the list of built-in scripts and exit immediately.

Optionally, write an integer weight after `@` to adjust the probability of selecting this script versus other ones. The default weight is 1. See below for details.

`-c clients`

`--client=clients`

Number of clients simulated, that is, number of concurrent database sessions. Default is 1.

`-C`

`--connect`

Establish a new connection for each transaction, rather than doing it just once per client session. This is useful to measure the connection overhead.

`-D varname=value`

`--define=varname=value`

Define a variable for use by a custom script (see below). Multiple `-D` options are allowed.

`-f filename[@weight]`

`--file=filename[@weight]`

Add a transaction script read from *filename* to the list of scripts to be executed.

Optionally, write an integer weight after `@` to adjust the probability of selecting this script versus other ones. The default weight is 1. (To use a script file name that includes an `@` character, append a weight so that there is no ambiguity, for example `filen@me@1`.) See below for details.

`-j threads`

`--jobs=threads`

Number of worker threads within `pgbench`. Using more than one thread can be helpful on multi-CPU machines. Clients are distributed as evenly as possible among available threads. Default is 1.

-l
--log

Write information about each transaction to a log file. See below for details.

-L *limit*
--latency-limit=*limit*

Transactions that last more than *limit* milliseconds are counted and reported separately, as *late*.

When throttling is used (`--rate=...`), transactions that lag behind schedule by more than *limit* ms, and thus have no hope of meeting the latency limit, are not sent to the server at all. They are counted and reported separately as *skipped*.

When the `--max-tries` option is used, a transaction which fails due to a serialization anomaly or from a deadlock will not be retried if the total time of all its tries is greater than *limit* ms. To limit only the time of tries and not their number, use `--max-tries=0`. By default, the option `--max-tries` is set to 1 and transactions with serialization/deadlock errors are not retried. See [Failures and Serialization/Deadlock Retries](#) for more information about retrying such transactions.

-M *querymode*
--protocol=*querymode*

Protocol to use for submitting queries to the server:

- `simple`: use simple query protocol.
- `extended`: use extended query protocol.
- `prepared`: use extended query protocol with prepared statements.

In the `prepared` mode, `pgbench` reuses the parse analysis result starting from the second query iteration, so `pgbench` runs faster than in other modes.

The default is simple query protocol. (See [Chapter 56](#) for more information.)

-n
--no-vacuum

Perform no vacuuming before running the test. This option is *necessary* if you are running a custom test scenario that does not include the standard tables `pgbench_accounts`, `pgbench_branches`, `pgbench_history`, and `pgbench_tellers`.

-N
--skip-some-updates

Run built-in simple-update script. Shorthand for `-b simple-update`.

-P *sec*
--progress=*sec*

Show progress report every *sec* seconds. The report includes the time since the beginning of the run, the TPS since the last report, and the transaction latency average, standard deviation, and the number of failed transactions since the last report. Under throttling (`-R`), the latency is computed with respect to the transaction scheduled start time, not the actual transaction beginning time, thus it also includes the average schedule lag time. When `--max-tries` is used to enable transaction retries after serialization/deadlock errors, the report includes the number of retried transactions and the sum of all retries.

-r
--report-per-command

Report the following statistics for each command after the benchmark finishes: the average per-statement latency (execution time from the perspective of the client), the number of failures, and the number of retries after serialization or deadlock errors in this command. The report displays retry statistics only if the `--max-tries` option is not equal to 1.

`-R rate`
`--rate=rate`

Execute transactions targeting the specified rate instead of running as fast as possible (the default). The rate is given in transactions per second. If the targeted rate is above the maximum possible rate, the rate limit won't impact the results.

The rate is targeted by starting transactions along a Poisson-distributed schedule time line. The expected start time schedule moves forward based on when the client first started, not when the previous transaction ended. That approach means that when transactions go past their original scheduled end time, it is possible for later ones to catch up again.

When throttling is active, the transaction latency reported at the end of the run is calculated from the scheduled start times, so it includes the time each transaction had to wait for the previous transaction to finish. The wait time is called the schedule lag time, and its average and maximum are also reported separately. The transaction latency with respect to the actual transaction start time, i.e., the time spent executing the transaction in the database, can be computed by subtracting the schedule lag time from the reported latency.

If `--latency-limit` is used together with `--rate`, a transaction can lag behind so much that it is already over the latency limit when the previous transaction ends, because the latency is calculated from the scheduled start time. Such transactions are not sent to the server, but are skipped altogether and counted separately.

A high schedule lag time is an indication that the system cannot process transactions at the specified rate, with the chosen number of clients and threads. When the average transaction execution time is longer than the scheduled interval between each transaction, each successive transaction will fall further behind, and the schedule lag time will keep increasing the longer the test run is. When that happens, you will have to reduce the specified transaction rate.

`-s scale_factor`
`--scale=scale_factor`

Report the specified scale factor in pgbench's output. With the built-in tests, this is not necessary; the correct scale factor will be detected by counting the number of rows in the `pgbench_branches` table. However, when testing only custom benchmarks (`-f` option), the scale factor will be reported as 1 unless this option is used.

`-S`
`--select-only`

Run built-in select-only script. Shorthand for `-b select-only`.

`-t transactions`
`--transactions=transactions`

Number of transactions each client runs. Default is 10.

`-T seconds`
`--time=seconds`

Run the test for this many seconds, rather than a fixed number of transactions per client. `-t` and `-T` are mutually exclusive.

`-v`
`--vacuum-all`

Vacuum all four standard tables before running the test. With neither `-n` nor `-v`, pgbench will vacuum the `pgbench_tellers` and `pgbench_branches` tables, and will truncate `pgbench_history`.

`--aggregate-interval=seconds`

Length of aggregation interval (in seconds). May be used only with `-l` option. With this option, the log contains per-interval summary data, as described below.

`--exit-on-abort`

Exit immediately when any client is aborted due to some error. Without this option, even when a client is aborted, other clients could continue their run as specified by `-t` or `-T` option, and `pgbench` will print an incomplete results in this case.

Note that serialization failures or deadlock failures do not abort the client, so they are not affected by this option. See [Failures and Serialization/Deadlock Retries](#) for more information.

`--failures-detailed`

Report failures in per-transaction and aggregation logs, as well as in the main and per-script reports, grouped by the following types:

- serialization failures;
- deadlock failures;

See [Failures and Serialization/Deadlock Retries](#) for more information.

`--log-prefix=prefix`

Set the filename prefix for the log files created by `--log`. The default is `pgbench_log`.

`--max-tries=number_of_tries`

Enable retries for transactions with serialization/deadlock errors and set the maximum number of these tries. This option can be combined with the `--latency-limit` option which limits the total time of all transaction tries; moreover, you cannot use an unlimited number of tries (`--max-tries=0`) without `--latency-limit` or `--time`. The default value is 1 and transactions with serialization/deadlock errors are not retried. See [Failures and Serialization/Deadlock Retries](#) for more information about retrying such transactions.

`--progress-timestamp`

When showing progress (option `-P`), use a timestamp (Unix epoch) instead of the number of seconds since the beginning of the run. The unit is in seconds, with millisecond precision after the dot. This helps compare logs generated by various tools.

`--random-seed=seed`

Set random generator seed. Seeds the system random number generator, which then produces a sequence of initial generator states, one for each thread. Values for *seed* may be: `time` (the default, the seed is based on the current time), `rand` (use a strong random source, failing if none is available), or an unsigned decimal integer value. The random generator is invoked explicitly from a `pgbench` script (`random...` functions) or implicitly (for instance option `--rate` uses it to schedule transactions). When explicitly set, the value used for seeding is shown on the terminal. Any value allowed for *seed* may also be provided through the environment variable `PGBENCH_RANDOM_SEED`. To ensure that the provided seed impacts all possible uses, put this option first or use the environment variable.

Setting the seed explicitly allows to reproduce a `pgbench` run exactly, as far as random numbers are concerned. As the random state is managed per thread, this means the exact same `pgbench` run for an identical invocation if there is one client per thread and there are no external or data dependencies. From a statistical viewpoint reproducing runs exactly is a bad idea because it can hide the performance variability or improve performance unduly, e.g., by hitting the same pages as a previous run. However, it may also be of great help for debugging, for instance re-running a tricky case which leads to an error. Use wisely.

`--sampling-rate=rate`

Sampling rate, used when writing data into the log, to reduce the amount of log generated. If this option is given, only the specified fraction of transactions are logged. 1.0 means all transactions will be logged, 0.05 means only 5% of the transactions will be logged.

Remember to take the sampling rate into account when processing the log file. For example, when computing TPS values, you need to multiply the numbers accordingly (e.g., with 0.01 sample rate, you'll only get 1/100 of the actual TPS).

`--show-script=scriptname`

Show the actual code of builtin script *scriptname* on stderr, and exit immediately.

`--verbose-errors`

Print messages about all errors and failures (errors without retrying) including which limit for retries was exceeded and how far it was exceeded for the serialization/deadlock failures. (Note that in this case the output can be significantly increased.) See [Failures and Serialization/Deadlock Retries](#) for more information.

Common Options

pgbench also accepts the following common command-line arguments for connection parameters and other common settings:

`--debug`

Print debugging output.

`-h hostname`

`--host=hostname`

The database server's host name

`-p port`

`--port=port`

The database server's port number

`-U login`

`--username=login`

The user name to connect as

`-V`

`--version`

Print the pgbench version and exit.

`-?`

`--help`

Show help about pgbench command line arguments, and exit.

Exit Status

A successful run will exit with status 0. Exit status 1 indicates static problems such as invalid command-line options or internal errors which are supposed to never occur. Early errors that occur when starting benchmark such as initial connection failures also exit with status 1. Errors during the run such as database errors or problems in the script will result in exit status 2. In the latter case, pgbench will print partial results if `--exit-on-abort` option is not specified.

Environment

PGDATABASE

PGHOST

PGPORT

PGUSER

Default connection parameters.

This utility, like most other Postgres Pro utilities, uses the environment variables supported by libpq (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

What Is the “Transaction” Actually Performed in pgbench?

pgbench executes test scripts chosen randomly from a specified list. The scripts may include built-in scripts specified with `-b` and user-provided scripts specified with `-f`. Each script may be given a relative weight specified after an `@` so as to change its selection probability. The default weight is 1. Scripts with a weight of 0 are ignored.

The default built-in transaction script (also invoked with `-b tpcb-like`) issues seven commands per transaction over randomly chosen `aid`, `tid`, `bid` and `delta`. The scenario is inspired by the TPC-B benchmark, but is not actually TPC-B, hence the name.

1. `BEGIN;`
2. `UPDATE pgbench_accounts SET abalance = abalance + :delta WHERE aid = :aid;`
3. `SELECT abalance FROM pgbench_accounts WHERE aid = :aid;`
4. `UPDATE pgbench_tellers SET tbalance = tbalance + :delta WHERE tid = :tid;`
5. `UPDATE pgbench_branches SET bbalance = bbalance + :delta WHERE bid = :bid;`
6. `INSERT INTO pgbench_history (tid, bid, aid, delta, mtime) VALUES (:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);`
7. `END;`

If you select the `simple-update` built-in (also `-N`), steps 4 and 5 aren't included in the transaction. This will avoid update contention on these tables, but it makes the test case even less like TPC-B.

If you select the `select-only` built-in (also `-S`), only the `SELECT` is issued.

Custom Scripts

pgbench has support for running custom benchmark scenarios by replacing the default transaction script (described above) with a transaction script read from a file (`-f` option). In this case a “transaction” counts as one execution of a script file.

A script file contains one or more SQL commands terminated by semicolons. Empty lines and lines beginning with `--` are ignored. Script files can also contain “meta commands”, which are interpreted by pgbench itself, as described below.

Note

Before Postgres Pro 9.6, SQL commands in script files were terminated by newlines, and so they could not be continued across lines. Now a semicolon is *required* to separate consecutive SQL commands (though an SQL command does not need one if it is followed by a meta command). If you need to create a script file that works with both old and new versions of pgbench, be sure to write each SQL command on a single line ending with a semicolon.

It is assumed that pgbench scripts do not contain incomplete blocks of SQL transactions. If at runtime the client reaches the end of the script without completing the last transaction block, it will be aborted.

There is a simple variable-substitution facility for script files. Variable names must consist of letters (including non-Latin letters), digits, and underscores, with the first character not being a digit. Variables can be set by the command-line `-D` option, explained above, or by the meta commands explained below. In addition to any variables preset by `-D` command-line options, there are a few variables that are preset automatically, listed in [Table 312](#). A value specified for these variables using `-D` takes precedence over the automatic presets. Once set, a variable's value can be inserted into an SQL command by writing `:variablename`. When running more than one client session, each session has its own set of variables. `pgbench` supports up to 255 variable uses in one statement.

Table 312. `pgbench` Automatic Variables

Variable	Description
<code>client_id</code>	unique number identifying the client session (starts from zero)
<code>default_seed</code>	seed used in hash and pseudorandom permutation functions by default
<code>random_seed</code>	random generator seed (unless overwritten with <code>-D</code>)
<code>scale</code>	current scale factor

Script file meta commands begin with a backslash (`\`) and normally extend to the end of the line, although they can be continued to additional lines by writing backslash-return. Arguments to a meta command are separated by white space. These meta commands are supported:

```
\gset [prefix] \aset [prefix]
```

These commands may be used to end SQL queries, taking the place of the terminating semicolon (`;`).

When the `\gset` command is used, the preceding SQL query is expected to return one row, the columns of which are stored into variables named after column names, and prefixed with *prefix* if provided.

When the `\aset` command is used, all combined SQL queries (separated by `\;`) have their columns stored into variables named after column names, and prefixed with *prefix* if provided. If a query returns no row, no assignment is made and the variable can be tested for existence to detect this. If a query returns more than one row, the last value is kept.

`\gset` and `\aset` cannot be used in pipeline mode, since the query results are not yet available by the time the commands would need them.

The following example puts the final account balance from the first query into variable *abalance*, and fills variables *p_two* and *p_three* with integers from the third query. The result of the second query is discarded. The result of the two last combined queries are stored in variables *four* and *five*.

```
UPDATE pgbench_accounts
SET abalance = abalance + :delta
WHERE aid = :aid
RETURNING abalance \gset
-- compound of two queries
SELECT 1 \;
SELECT 2 AS two, 3 AS three \gset p_
SELECT 4 AS four \; SELECT 5 AS five \aset
```

```
\if expression
\elif expression
\else
\endif
```

This group of commands implements nestable conditional blocks, similarly to `psql`'s `\if expression`. Conditional expressions are identical to those with `\set`, with non-zero values interpreted as true.

```
\set varname expression
```

Sets variable *varname* to a value calculated from *expression*. The expression may contain the `NULL` constant, Boolean constants `TRUE` and `FALSE`, integer constants such as `5432`, double constants such as `3.14159`, references to variables `:variablename`, [operators](#) with their usual SQL precedence and associativity, [function calls](#), SQL [CASE generic conditional expressions](#) and parentheses.

Functions and most operators return `NULL` on `NULL` input.

For conditional purposes, non zero numerical values are `TRUE`, zero numerical values and `NULL` are `FALSE`.

Too large or small integer and double constants, as well as integer arithmetic operators (`+`, `-`, `*` and `/`) raise errors on overflows.

When no final `ELSE` clause is provided to a `CASE`, the default value is `NULL`.

Examples:

```
\set ntellers 10 * :scale
\set aid (1021 * random(1, 100000 * :scale)) % \
        (100000 * :scale) + 1
\set divx CASE WHEN :x <> 0 THEN :y/:x ELSE NULL END
```

```
\sleep number [ us | ms | s ]
```

Causes script execution to sleep for the specified duration in microseconds (`us`), milliseconds (`ms`) or seconds (`s`). If the unit is omitted then seconds are the default. *number* can be either an integer constant or a `:variablename` reference to a variable having an integer value.

Example:

```
\sleep 10 ms
```

```
\setshell varname command [ argument ... ]
```

Sets variable *varname* to the result of the shell command *command* with the given *argument(s)*. The command must return an integer value through its standard output.

command and each *argument* can be either a text constant or a `:variablename` reference to a variable. If you want to use an *argument* starting with a colon, write an additional colon at the beginning of *argument*.

Example:

```
\setshell variable_to_be_assigned command
literal_argument :variable ::literal_starting_with_colon
```

```
\shell command [ argument ... ]
```

Same as `\setshell`, but the result of the command is discarded.

Example:

```
\shell command literal_argument :variable ::literal_starting_with_colon
```

```
\startpipeline
```

```
\syncpipeline
```

```
\endpipeline
```

This group of commands implements pipelining of SQL statements. A pipeline must begin with a `\startpipeline` and end with an `\endpipeline`. In between there may be any number of `\syncpipeline` commands, which sends a [sync message](#) without ending the ongoing pipeline and

flushing the send buffer. In pipeline mode, statements are sent to the server without waiting for the results of previous statements. See [Section 35.5](#) for more details. Pipeline mode requires the use of extended query protocol.

Built-in Operators

The arithmetic, bitwise, comparison and logical operators listed in [Table 313](#) are built into pgbench and may be used in expressions appearing in `\set`. The operators are listed in increasing precedence order. Except as noted, operators taking two numeric inputs will produce a double value if either input is double, otherwise they produce an integer result.

Table 313. pgbench Operators

Operator Description Example(s)
<i>boolean OR boolean → boolean</i> Logical OR <code>5 or 0 → TRUE</code>
<i>boolean AND boolean → boolean</i> Logical AND <code>3 and 0 → FALSE</code>
<i>NOT boolean → boolean</i> Logical NOT <code>not false → TRUE</code>
<i>boolean IS [NOT] (NULL TRUE FALSE) → boolean</i> Boolean value tests <code>1 is null → FALSE</code>
<i>value ISNULL NOTNULL → boolean</i> Nullness tests <code>1 notnull → TRUE</code>
<i>number = number → boolean</i> Equal <code>5 = 4 → FALSE</code>
<i>number <> number → boolean</i> Not equal <code>5 <> 4 → TRUE</code>
<i>number != number → boolean</i> Not equal <code>5 != 5 → FALSE</code>
<i>number < number → boolean</i> Less than <code>5 < 4 → FALSE</code>
<i>number <= number → boolean</i> Less than or equal to <code>5 <= 4 → FALSE</code>
<i>number > number → boolean</i> Greater than <code>5 > 4 → TRUE</code>

Operator	Description	Example(s)
$number \geq number \rightarrow boolean$	Greater than or equal to	$5 \geq 4 \rightarrow TRUE$
$integer \mid integer \rightarrow integer$	Bitwise OR	$1 \mid 2 \rightarrow 3$
$integer \# integer \rightarrow integer$	Bitwise XOR	$1 \# 3 \rightarrow 2$
$integer \& integer \rightarrow integer$	Bitwise AND	$1 \& 3 \rightarrow 1$
$\sim integer \rightarrow integer$	Bitwise NOT	$\sim 1 \rightarrow -2$
$integer \ll integer \rightarrow integer$	Bitwise shift left	$1 \ll 2 \rightarrow 4$
$integer \gg integer \rightarrow integer$	Bitwise shift right	$8 \gg 2 \rightarrow 2$
$number + number \rightarrow number$	Addition	$5 + 4 \rightarrow 9$
$number - number \rightarrow number$	Subtraction	$3 - 2.0 \rightarrow 1.0$
$number * number \rightarrow number$	Multiplication	$5 * 4 \rightarrow 20$
$number / number \rightarrow number$	Division (truncates the result towards zero if both inputs are integers)	$5 / 3 \rightarrow 1$
$integer \% integer \rightarrow integer$	Modulo (remainder)	$3 \% 2 \rightarrow 1$
$- number \rightarrow number$	Negation	$- 2.0 \rightarrow -2.0$

Built-In Functions

The functions listed in [Table 314](#) are built into pgbench and may be used in expressions appearing in `\set`.

Table 314. pgbench Functions

Function	Description	Example(s)
<code>abs (number)</code>	→ same type as input Absolute value	<code>abs (-17) → 17</code>
<code>debug (number)</code>	→ same type as input Prints the argument to stderr, and returns the argument.	<code>debug (5432.1) → 5432.1</code>
<code>double (number)</code>	→ double Casts to double.	<code>double (5432) → 5432.0</code>
<code>exp (number)</code>	→ double Exponential (e raised to the given power)	<code>exp (1.0) → 2.718281828459045</code>
<code>greatest (number [, ...])</code>	→ double if any argument is double, else integer Selects the largest value among the arguments.	<code>greatest (5, 4, 3, 2) → 5</code>
<code>hash (value [, seed])</code>	→ integer This is an alias for <code>hash_murmur2</code> .	<code>hash (10, 5432) → -5817877081768721676</code>
<code>hash_fnv1a (value [, seed])</code>	→ integer Computes FNV-1a hash .	<code>hash_fnv1a (10, 5432) → -7793829335365542153</code>
<code>hash_murmur2 (value [, seed])</code>	→ integer Computes MurmurHash2 hash .	<code>hash_murmur2 (10, 5432) → -5817877081768721676</code>
<code>int (number)</code>	→ integer Casts to integer.	<code>int (5.4 + 3.8) → 9</code>
<code>least (number [, ...])</code>	→ double if any argument is double, else integer Selects the smallest value among the arguments.	<code>least (5, 4, 3, 2.1) → 2.1</code>
<code>ln (number)</code>	→ double Natural logarithm	<code>ln (2.718281828459045) → 1.0</code>
<code>mod (integer, integer)</code>	→ integer Modulo (remainder)	<code>mod (54, 32) → 22</code>
<code>permute (i, size [, seed])</code>	→ integer Permuted value of <i>i</i> , in the range <code>[0, size)</code> . This is the new position of <i>i</i> (modulo <i>size</i>) in a pseudorandom permutation of the integers <code>0...size-1</code> , parameterized by <i>seed</i> , see below.	<code>permute (0, 4) → an integer between 0 and 3</code>
<code>pi ()</code>	→ double	

Function	Description	Example(s)
	Approximate value of π	<code>pi()</code> \rightarrow 3.14159265358979323846
	<code>pow(x, y)</code> \rightarrow double <code>power(x, y)</code> \rightarrow double x raised to the power of y	<code>pow(2.0, 10)</code> \rightarrow 1024.0
	<code>random(lb, ub)</code> \rightarrow integer Computes a uniformly-distributed random integer in $[lb, ub]$.	<code>random(1, 10)</code> \rightarrow an integer between 1 and 10
	<code>random_exponential(lb, ub, parameter)</code> \rightarrow integer Computes an exponentially-distributed random integer in $[lb, ub]$, see below.	<code>random_exponential(1, 10, 3.0)</code> \rightarrow an integer between 1 and 10
	<code>random_gaussian(lb, ub, parameter)</code> \rightarrow integer Computes a Gaussian-distributed random integer in $[lb, ub]$, see below.	<code>random_gaussian(1, 10, 2.5)</code> \rightarrow an integer between 1 and 10
	<code>random_zipfian(lb, ub, parameter)</code> \rightarrow integer Computes a Zipfian-distributed random integer in $[lb, ub]$, see below.	<code>random_zipfian(1, 10, 1.5)</code> \rightarrow an integer between 1 and 10
	<code>sqrt(number)</code> \rightarrow double Square root	<code>sqrt(2.0)</code> \rightarrow 1.414213562

The `random` function generates values using a uniform distribution, that is all the values are drawn within the specified range with equal probability. The `random_exponential`, `random_gaussian` and `random_zipfian` functions require an additional double parameter which determines the precise shape of the distribution.

- For an exponential distribution, *parameter* controls the distribution by truncating a quickly-decreasing exponential distribution at *parameter*, and then projecting onto integers between the bounds. To be precise, with

$$f(x) = \exp(-\text{parameter} * (x - \min) / (\max - \min + 1)) / (1 - \exp(-\text{parameter}))$$

Then value *i* between *min* and *max* inclusive is drawn with probability: $f(i) - f(i + 1)$.

Intuitively, the larger the *parameter*, the more frequently values close to *min* are accessed, and the less frequently values close to *max* are accessed. The closer to 0 *parameter* is, the flatter (more uniform) the access distribution. A crude approximation of the distribution is that the most frequent 1% values in the range, close to *min*, are drawn *parameter*% of the time. The *parameter* value must be strictly positive.

- For a Gaussian distribution, the interval is mapped onto a standard normal distribution (the classical bell-shaped Gaussian curve) truncated at $-\text{parameter}$ on the left and $+\text{parameter}$ on the right. Values in the middle of the interval are more likely to be drawn. To be precise, if $\text{PHI}(x)$ is the cumulative distribution function of the standard normal distribution, with mean *mu* defined as $(\max + \min) / 2.0$, with

$$f(x) = \text{PHI}(2.0 * \text{parameter} * (x - \mu) / (\max - \min + 1)) / (2.0 * \text{PHI}(\text{parameter}) - 1)$$

then value i between min and max inclusive is drawn with probability: $f(i + 0.5) - f(i - 0.5)$. Intuitively, the larger the *parameter*, the more frequently values close to the middle of the interval are drawn, and the less frequently values close to the *min* and *max* bounds. About 67% of values are drawn from the middle $1.0 / parameter$, that is a relative $0.5 / parameter$ around the mean, and 95% in the middle $2.0 / parameter$, that is a relative $1.0 / parameter$ around the mean; for instance, if *parameter* is 4.0, 67% of values are drawn from the middle quarter ($1.0 / 4.0$) of the interval (i.e., from $3.0 / 8.0$ to $5.0 / 8.0$) and 95% from the middle half ($2.0 / 4.0$) of the interval (second and third quartiles). The minimum allowed *parameter* value is 2.0.

- `random_zipfian` generates a bounded Zipfian distribution. *parameter* defines how skewed the distribution is. The larger the *parameter*, the more frequently values closer to the beginning of the interval are drawn. The distribution is such that, assuming the range starts from 1, the ratio of the probability of drawing k versus drawing $k+1$ is $((k+1)/k)**parameter$. For example, `random_zipfian(1, ..., 2.5)` produces the value 1 about $(2/1)**2.5 = 5.66$ times more frequently than 2, which itself is produced $(3/2)**2.5 = 2.76$ times more frequently than 3, and so on.

pgbench's implementation is based on "Non-Uniform Random Variate Generation", Luc Devroye, p. 550-551, Springer 1986. Due to limitations of that algorithm, the *parameter* value is restricted to the range [1.001, 1000].

Note

When designing a benchmark which selects rows non-uniformly, be aware that the rows chosen may be correlated with other data such as IDs from a sequence or the physical row ordering, which may skew performance measurements.

To avoid this, you may wish to use the `permute` function, or some other additional step with similar effect, to shuffle the selected rows and remove such correlations.

Hash functions `hash`, `hash_murmur2` and `hash_fnv1a` accept an input value and an optional seed parameter. In case the seed isn't provided the value of `:default_seed` is used, which is initialized randomly unless set by the command-line `-D` option.

`permute` accepts an input value, a size, and an optional seed parameter. It generates a pseudorandom permutation of integers in the range $[0, size)$, and returns the index of the input value in the permuted values. The permutation chosen is parameterized by the seed, which defaults to `:default_seed`, if not specified. Unlike the hash functions, `permute` ensures that there are no collisions or holes in the output values. Input values outside the interval are interpreted modulo the size. The function raises an error if the size is not positive. `permute` can be used to scatter the distribution of non-uniform random functions such as `random_zipfian` or `random_exponential` so that values drawn more often are not trivially correlated. For instance, the following pgbench script simulates a possible real world workload typical for social media and blogging platforms where a few accounts generate excessive load:

```
\set size 1000000
\set r random_zipfian(1, :size, 1.07)
\set k 1 + permute(:r, :size)
```

In some cases several distinct distributions are needed which don't correlate with each other and this is when the optional seed parameter comes in handy:

```
\set k1 1 + permute(:r, :size, :default_seed + 123)
\set k2 1 + permute(:r, :size, :default_seed + 321)
```

A similar behavior can also be approximated with `hash`:

```
\set size 1000000
\set r random_zipfian(1, 100 * :size, 1.07)
\set k 1 + abs(hash(:r)) % :size
```

However, since `hash` generates collisions, some values will not be reachable and others will be more frequent than expected from the original distribution.

As an example, the full definition of the built-in TPC-B-like transaction is:

```
\set aid random(1, 100000 * :scale)
\set bid random(1, 1 * :scale)
\set tid random(1, 10 * :scale)
\set delta random(-5000, 5000)
BEGIN;
UPDATE pgbench_accounts SET abalance = abalance + :delta WHERE aid = :aid;
SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
UPDATE pgbench_tellers SET tbalance = tbalance + :delta WHERE tid = :tid;
UPDATE pgbench_branches SET bbalance = bbalance + :delta WHERE bid = :bid;
INSERT INTO pgbench_history (tid, bid, aid, delta, mtime) VALUES
(:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
END;
```

This script allows each iteration of the transaction to reference different, randomly-chosen rows. (This example also shows why it's important for each client session to have its own variables — otherwise they'd not be independently touching different rows.)

Per-Transaction Logging

With the `-l` option (but without the `--aggregate-interval` option), `pgbench` writes information about each transaction to a log file. The log file will be named *prefix.nnn*, where *prefix* defaults to `pgbench_log`, and *nnn* is the PID of the `pgbench` process. The prefix can be changed by using the `--log-prefix` option. If the `-j` option is 2 or higher, so that there are multiple worker threads, each will have its own log file. The first worker will use the same name for its log file as in the standard single worker case. The additional log files for the other workers will be named *prefix.nnn.mmm*, where *mmm* is a sequential number for each worker starting with 1.

Each line in a log file describes one transaction. It contains the following space-separated fields:

client_id

identifies the client session that ran the transaction

transaction_no

counts how many transactions have been run by that session

time

transaction's elapsed time, in microseconds

script_no

identifies the script file that was used for the transaction (useful when multiple scripts are specified with `-f` or `-b`)

time_epoch

transaction's completion time, as a Unix-epoch time stamp

time_us

fractional-second part of transaction's completion time, in microseconds

schedule_lag

transaction start delay, that is the difference between the transaction's scheduled start time and the time it actually started, in microseconds (present only if `--rate` is specified)

retries

count of retries after serialization or deadlock errors during the transaction (present only if `--max-tries` is not equal to one)

When both `--rate` and `--latency-limit` are used, the *time* for a skipped transaction will be reported as `skipped`. If the transaction ends with a failure, its *time* will be reported as `failed`. If you use the `--failures-detailed` option, the *time* of the failed transaction will be reported as `serialization` or `deadlock` depending on the type of failure (see [Failures and Serialization/Deadlock Retries](#) for more information).

Here is a snippet of a log file generated in a single-client run:

```
0 199 2241 0 1175850568 995598
0 200 2465 0 1175850568 998079
0 201 2513 0 1175850569 608
0 202 2038 0 1175850569 2663
```

Another example with `--rate=100` and `--latency-limit=5` (note the additional *schedule_lag* column):

```
0 81 4621 0 1412881037 912698 3005
0 82 6173 0 1412881037 914578 4304
0 83 skipped 0 1412881037 914578 5217
0 83 skipped 0 1412881037 914578 5099
0 83 4722 0 1412881037 916203 3108
0 84 4142 0 1412881037 918023 2333
0 85 2465 0 1412881037 919759 740
```

In this example, transaction 82 was late, because its latency (6.173 ms) was over the 5 ms limit. The next two transactions were skipped, because they were already late before they were even started.

The following example shows a snippet of a log file with failures and retries, with the maximum number of tries set to 10 (note the additional *retries* column):

```
3 0 47423 0 1499414498 34501 3
3 1 8333 0 1499414498 42848 0
3 2 8358 0 1499414498 51219 0
4 0 72345 0 1499414498 59433 6
1 3 41718 0 1499414498 67879 4
1 4 8416 0 1499414498 76311 0
3 3 33235 0 1499414498 84469 3
0 0 failed 0 1499414498 84905 9
2 0 failed 0 1499414498 86248 9
3 4 8307 0 1499414498 92788 0
```

If the `--failures-detailed` option is used, the type of failure is reported in the *time* like this:

```
3 0 47423 0 1499414498 34501 3
3 1 8333 0 1499414498 42848 0
3 2 8358 0 1499414498 51219 0
4 0 72345 0 1499414498 59433 6
1 3 41718 0 1499414498 67879 4
1 4 8416 0 1499414498 76311 0
3 3 33235 0 1499414498 84469 3
0 0 serialization 0 1499414498 84905 9
2 0 serialization 0 1499414498 86248 9
3 4 8307 0 1499414498 92788 0
```

When running a long test on hardware that can handle a lot of transactions, the log files can become very large. The `--sampling-rate` option can be used to log only a random sample of transactions.

Aggregated Logging

With the `--aggregate-interval` option, a different format is used for the log files. Each log line describes one aggregation interval. It contains the following space-separated fields:

interval_start

start time of the interval, as a Unix-epoch time stamp

num_transactions

number of transactions within the interval

sum_latency

sum of transaction latencies

sum_latency_2

sum of squares of transaction latencies

min_latency

minimum transaction latency

max_latency

maximum transaction latency

sum_lag

sum of transaction start delays (zero unless `--rate` is specified)

sum_lag_2

sum of squares of transaction start delays (zero unless `--rate` is specified)

min_lag

minimum transaction start delay (zero unless `--rate` is specified)

max_lag

maximum transaction start delay (zero unless `--rate` is specified)

skipped

number of transactions skipped because they would have started too late (zero unless `--rate` and `--latency-limit` are specified)

retried

number of retried transactions (zero unless `--max-tries` is not equal to one)

retries

number of retries after serialization or deadlock errors (zero unless `--max-tries` is not equal to one)

serialization_failures

number of transactions that got a serialization error and were not retried afterwards (zero unless `--failures-detailed` is specified)

deadlock_failures

number of transactions that got a deadlock error and were not retried afterwards (zero unless `--failures-detailed` is specified)

Here is some example output generated with this option:

```
pgbench --aggregate-interval=10 --time=20 --client=10 --log --rate=1000 --latency-limit=10 --failures-detailed --max-tries=10 test
```

```
1650260552 5178 26171317 177284491527 1136 44462 2647617 7321113867 0 9866 64 7564
28340 4148 0
1650260562 4808 25573984 220121792172 1171 62083 3037380 9666800914 0 9998 598 7392
26621 4527 0
```

Notice that while the plain (unaggregated) log format shows which script was used for each transaction, the aggregated format does not. Therefore if you need per-script data, you need to aggregate the data on your own.

Per-Statement Report

With the `-r` option, pgbench collects the following statistics for each statement:

- `latency` — elapsed transaction time for each statement. pgbench reports an average value of all successful runs of the statement.
- The number of failures in this statement. See [Failures and Serialization/Deadlock Retries](#) for more information.
- The number of retries after a serialization or a deadlock error in this statement. See [Failures and Serialization/Deadlock Retries](#) for more information.

The report displays retry statistics only if the `--max-tries` option is not equal to 1.

All values are computed for each statement executed by every client and are reported after the benchmark has finished.

For the default script, the output will look similar to this:

```
starting vacuum...end.
transaction type: <builtin: TPC-B (sort of)>
scaling factor: 1
query mode: simple
number of clients: 10
number of threads: 1
maximum number of tries: 1
number of transactions per client: 1000
number of transactions actually processed: 10000/10000
number of failed transactions: 0 (0.000%)
number of transactions above the 50.0 ms latency limit: 1311/10000 (13.110 %)
latency average = 28.488 ms
latency stddev = 21.009 ms
initial connection time = 69.068 ms
tps = 346.224794 (without initial connection time)
statement latencies in milliseconds and failures:
  0.012  0  \set aid random(1, 100000 * :scale)
  0.002  0  \set bid random(1, 1 * :scale)
  0.002  0  \set tid random(1, 10 * :scale)
  0.002  0  \set delta random(-5000, 5000)
  0.319  0  BEGIN;
  0.834  0  UPDATE pgbench_accounts SET abalance = abalance + :delta WHERE aid = :aid;
  0.641  0  SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
 11.126  0  UPDATE pgbench_tellers SET tbalance = tbalance + :delta WHERE tid = :tid;
 12.961  0  UPDATE pgbench_branches SET bbalance = bbalance + :delta WHERE bid = :bid;
  0.634  0  INSERT INTO pgbench_history (tid, bid, aid, delta, mtime) VALUES
(:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
 1.957  0  END;
```

Another example of output for the default script using serializable default transaction isolation level (PGOPTIONS='-c default_transaction_isolation=serializable' pgbench ...):

```
starting vacuum...end.
transaction type: <builtin: TPC-B (sort of)>
scaling factor: 1
query mode: simple
number of clients: 10
number of threads: 1
```

```
maximum number of tries: 10
number of transactions per client: 1000
number of transactions actually processed: 6317/10000
number of failed transactions: 3683 (36.830%)
number of transactions retried: 7667 (76.670%)
total number of retries: 45339
number of transactions above the 50.0 ms latency limit: 106/6317 (1.678 %)
latency average = 17.016 ms
latency stddev = 13.283 ms
initial connection time = 45.017 ms
tps = 186.792667 (without initial connection time)
statement latencies in milliseconds, failures and retries:
 0.006      0      0  \set aid random(1, 100000 * :scale)
 0.001      0      0  \set bid random(1, 1 * :scale)
 0.001      0      0  \set tid random(1, 10 * :scale)
 0.001      0      0  \set delta random(-5000, 5000)
 0.385      0      0  BEGIN;
 0.773      0      1  UPDATE pgbench_accounts SET abalance = abalance + :delta WHERE
aid = :aid;
 0.624      0      0  SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
 1.098     320     3762  UPDATE pgbench_tellers SET tbalance = tbalance + :delta WHERE tid
= :tid;
 0.582    3363    41576  UPDATE pgbench_branches SET bbalance = bbalance + :delta WHERE
bid = :bid;
 0.465      0      0  INSERT INTO pgbench_history (tid, bid, aid, delta, mtime) VALUES
(:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
 1.933      0      0  END;
```

If multiple script files are specified, all statistics are reported separately for each script file.

Note that collecting the additional timing information needed for per-statement latency computation adds some overhead. This will slow average execution speed and lower the computed TPS. The amount of slowdown varies significantly depending on platform and hardware. Comparing average TPS values with and without latency reporting enabled is a good way to measure if the timing overhead is significant.

Failures and Serialization/Deadlock Retries

When executing pgbench, there are three main types of errors:

- Errors of the main program. They are the most serious and always result in an immediate exit from pgbench with the corresponding error message. They include:
 - errors at the beginning of pgbench (e.g. an invalid option value);
 - errors in the initialization mode (e.g. the query to create tables for built-in scripts fails);
 - errors before starting threads (e.g. could not connect to the database server, syntax error in the meta command, thread creation failure);
 - internal pgbench errors (which are supposed to never occur...).
- Errors when the thread manages its clients (e.g. the client could not start a connection to the database server / the socket for connecting the client to the database server has become invalid). In such cases all clients of this thread stop while other threads continue to work. However, `--exit-on-abort` is specified, all of the threads stop immediately in this case.
- Direct client errors. They lead to immediate exit from pgbench with the corresponding error message in the case of an internal pgbench error (which are supposed to never occur...) or when `--exit-on-abort` is specified. Otherwise in the worst case they only lead to the abortion of the failed client while other clients continue their run (but some client errors are handled without an abortion of the client and reported separately, see below). Later in this section it is assumed that the discussed errors are only the direct client errors and they are not internal pgbench errors.

A client's run is aborted in case of a serious error; for example, the connection with the database server was lost or the end of script was reached without completing the last transaction. In addition, if execution of an SQL or meta command fails for reasons other than serialization or deadlock errors, the client is aborted. Otherwise, if an SQL command fails with serialization or deadlock errors, the client is not aborted. In such cases, the current transaction is rolled back, which also includes setting the client variables as they were before the run of this transaction (it is assumed that one transaction script contains only one transaction; see [What Is the "Transaction" Actually Performed in pgbench?](#) for more information). Transactions with serialization or deadlock errors are repeated after rollbacks until they complete successfully or reach the maximum number of tries (specified by the `--max-tries` option) / the maximum time of retries (specified by the `--latency-limit` option) / the end of benchmark (specified by the `--time` option). If the last trial run fails, this transaction will be reported as failed but the client is not aborted and continues to work.

Note

Without specifying the `--max-tries` option, a transaction will never be retried after a serialization or deadlock error because its default value is 1. Use an unlimited number of tries (`--max-tries=0`) and the `--latency-limit` option to limit only the maximum time of tries. You can also use the `--time` option to limit the benchmark duration under an unlimited number of tries.

Be careful when repeating scripts that contain multiple transactions: the script is always retried completely, so successful transactions can be performed several times.

Be careful when repeating transactions with shell commands. Unlike the results of SQL commands, the results of shell commands are not rolled back, except for the variable value of the `\setshell` command.

The latency of a successful transaction includes the entire time of transaction execution with rollbacks and retries. The latency is measured only for successful transactions and commands but not for failed transactions or commands.

The main report contains the number of failed transactions. If the `--max-tries` option is not equal to 1, the main report also contains statistics related to retries: the total number of retried transactions and total number of retries. The per-script report inherits all these fields from the main report. The per-statement report displays retry statistics only if the `--max-tries` option is not equal to 1.

If you want to group failures by basic types in per-transaction and aggregation logs, as well as in the main and per-script reports, use the `--failures-detailed` option. If you also want to distinguish all errors and failures (errors without retrying) by type including which limit for retries was exceeded and how much it was exceeded by for the serialization/deadlock failures, use the `--verbose-errors` option.

Table Access Methods

You may specify the [Table Access Method](#) for the pgbench tables. The environment variable `PGOPTIONS` specifies database configuration options that are passed to Postgres Pro via the command line (See [Section 19.1.4](#)). For example, a hypothetical default Table Access Method for the tables that pgbench creates called `wuzza` can be specified with:

```
PGOPTIONS='-c default_table_access_method=wuzza'
```

Good Practices

It is very easy to use pgbench to produce completely meaningless numbers. Here are some guidelines to help you get useful results.

In the first place, *never* believe any test that runs for only a few seconds. Use the `-t` or `-T` option to make the run last at least a few minutes, so as to average out noise. In some cases you could need hours to get numbers that are reproducible. It's a good idea to try the test run a few times, to find out if your numbers are reproducible or not.

For the default TPC-B-like test scenario, the initialization scale factor (`-s`) should be at least as large as the largest number of clients you intend to test (`-c`); else you'll mostly be measuring update contention. There are only `-s` rows in the `pgbench_branches` table, and every transaction wants to update one of them, so `-c` values in excess of `-s` will undoubtedly result in lots of transactions blocked waiting for other transactions.

The default test scenario is also quite sensitive to how long it's been since the tables were initialized: accumulation of dead rows and dead space in the tables changes the results. To understand the results you must keep track of the total number of updates and when vacuuming happens. If `autovacuum` is enabled it can result in unpredictable changes in measured performance.

A limitation of `pgbench` is that it can itself become the bottleneck when trying to test a large number of client sessions. This can be alleviated by running `pgbench` on a different machine from the database server, although low network latency will be essential. It might even be useful to run several `pgbench` instances concurrently, on several client machines, against the same database server.

Security

If untrusted users have access to a database that has not adopted a [secure schema usage pattern](#), do not run `pgbench` in that database. `pgbench` uses unqualified names and does not manipulate the search path.

pg_combinebackup

`pg_combinebackup` — reconstruct a full backup from an incremental backup and dependent backups

Synopsis

`pg_combinebackup` [*option...*] [*backup_directory...*]

Description

`pg_combinebackup` is used to reconstruct a synthetic full backup from an incremental backup and the earlier backups upon which it depends.

Specify all of the required backups on the command line from oldest to newest. That is, the first backup directory should be the path to the full backup, and the last should be the path to the final incremental backup that you wish to restore. The reconstructed backup will be written to the output directory specified by the `-o` option.

`pg_combinebackup` will attempt to verify that the backups you specify form a legal backup chain from which a correct full backup can be reconstructed. However, it is not designed to help you keep track of which backups depend on which other backups. If you remove one or more of the previous backups upon which your incremental backup relies, you will not be able to restore it. Moreover, `pg_combinebackup` only attempts to verify that the backups have the correct relationship to each other, not that each individual backup is intact; for that, use [pg_verifybackup](#).

Since the output of `pg_combinebackup` is a synthetic full backup, it can be used as an input to a future invocation of `pg_combinebackup`. The synthetic full backup would be specified on the command line in lieu of the chain of backups from which it was reconstructed.

Options

`-d`
`--debug`

Print lots of debug logging output on `stderr`.

`-n`
`--dry-run`

The `-n/--dry-run` option instructs `pg_combinebackup` to figure out what would be done without actually creating the target directory or any output files. It is particularly useful in combination with `--debug`.

`-N`
`--no-sync`

By default, `pg_combinebackup` will wait for all files to be written safely to disk. This option causes `pg_combinebackup` to return without waiting, which is faster, but means that a subsequent operating system crash can leave the output backup corrupt. Generally, this option is useful for testing but should not be used when creating a production installation.

`-o outputdir`
`--output=outputdir`

Specifies the output directory to which the synthetic full backup should be written. Currently, this argument is required.

`-T olddir=newdir`
`--tablespace-mapping=olddir=newdir`

Relocates the tablespace in directory `olddir` to `newdir` during the backup. `olddir` is the absolute path of the tablespace as it exists in the final backup specified on the command line, and `newdir`

is the absolute path to use for the tablespace in the reconstructed backup. If either path needs to contain an equal sign (=), precede that with a backslash. This option can be specified multiple times for multiple tablespaces.

`--clone`

Use efficient file cloning (also known as “reflinks” on some systems) instead of copying files to the new data directory, which can result in near-instantaneous copying of the data files.

If a backup manifest is not available or does not contain checksum of the right type, file cloning will be used to copy the file, but the file will be also read block-by-block for the checksum calculation.

File cloning is only supported on some operating systems and file systems. If it is selected but not supported, the `pg_combinebackup` run will error. At present, it is supported on Linux (kernel 4.5 or later) with Btrfs and XFS (on file systems created with reflink support), and on macOS with APFS.

`--copy`

Perform regular file copy. This is the default. (See also `--copy-file-range` and `--clone`.)

`--copy-file-range`

Use the `copy_file_range` system call for efficient copying. On some file systems this gives results similar to `--clone`, sharing physical disk blocks, while on others it may still copy blocks, but do so via an optimized path. At present, it is supported on Linux and FreeBSD.

If a backup manifest is not available or does not contain checksum of the right type, `copy_file_range` will be used to copy the file, but the file will be also read block-by-block for the checksum calculation.

`--manifest-checksums=algorithm`

`pg_combinebackup` writes a backup manifest in the output directory. This option specifies the checksum algorithm that should be applied to each file included in the backup manifest. Currently, the available algorithms are `NONE`, `CRC32C`, `SHA224`, `SHA256`, `SHA384`, and `SHA512`. The default is `CRC32C`.

`--no-manifest`

Disables generation of a backup manifest. If this option is not specified, a backup manifest for the reconstructed backup will be written to the output directory.

`--sync-method=method`

When set to `fsync`, which is the default, `pg_combinebackup` will recursively open and synchronize all files in the backup directory. When the plain format is used, the search for files will follow symbolic links for the WAL directory and each configured tablespace.

On Linux, `syncfs` may be used instead to ask the operating system to synchronize the whole file system that contains the backup directory. When the plain format is used, `pg_combinebackup` will also synchronize the file systems that contain the WAL files and each tablespace. See [recovery_init_sync_method](#) for information about the caveats to be aware of when using `syncfs`.

This option has no effect when `--no-sync` is used.

`-V`

`--version`

Prints the `pg_combinebackup` version and exits.

`-?`

`--help`

Shows help about `pg_combinebackup` command line arguments, and exits.

Limitations

`pg_combinebackup` does not recompute page checksums when writing the output directory. Therefore, if any of the backups used for reconstruction were taken with checksums disabled, but the final backup was taken with checksums enabled, the resulting directory may contain pages with invalid checksums.

To avoid this problem, taking a new full backup after changing the checksum state of the cluster using [pg_checksums](#) is recommended. Otherwise, you can disable and then optionally reenable checksums on the directory produced by `pg_combinebackup` in order to correct the problem.

Environment

This utility, like most other PostgreSQL utilities, uses the environment variables supported by libpq (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

pg_config

pg_config — retrieve information about the installed version of Postgres Pro

Synopsis

```
pg_config [option...]
```

Description

The `pg_config` utility prints configuration parameters of the currently installed version of Postgres Pro. It is intended, for example, to be used by software packages that want to interface to Postgres Pro to facilitate finding the required header files and libraries.

Options

To use `pg_config`, supply one or more of the following options:

`--bindir`

Print the location of user executables. Use this, for example, to find the `psql` program. This is normally also the location where the `pg_config` program resides.

`--docdir`

Print the location of documentation files.

`--htmldir`

Print the location of HTML documentation files.

`--includedir`

Print the location of C header files of the client interfaces.

`--pkgincludedir`

Print the location of other C header files.

`--includedir-server`

Print the location of C header files for server programming.

`--libdir`

Print the location of object code libraries.

`--pkglibdir`

Print the location of dynamically loadable modules, or where the server would search for them. (Other architecture-dependent data files might also be installed in this directory.)

`--localedir`

Print the location of locale support files. (This will be an empty string if locale support was not configured when Postgres Pro was built.)

`--mandir`

Print the location of manual pages.

`--sharedir`

Print the location of architecture-independent support files.

`--sysconfdir`

Print the location of system-wide configuration files.

`--pgxs`

Print the location of extension makefiles.

`--configure`

Print the options that were given to the `configure` script when Postgres Pro was configured for building. This can be used to reproduce the identical configuration, or to find out with what options a binary package was built. (Note however that binary packages often contain vendor-specific custom patches.) See also the examples below.

`--cc`

Print the value of the `CC` variable that was used for building Postgres Pro. This shows the C compiler used.

`--cppflags`

Print the value of the `CPPFLAGS` variable that was used for building Postgres Pro. This shows C compiler switches needed at preprocessing time (typically, `-I` switches).

`--cflags`

Print the value of the `CFLAGS` variable that was used for building Postgres Pro. This shows C compiler switches.

`--cflags_sl`

Print the value of the `CFLAGS_SL` variable that was used for building Postgres Pro. This shows extra C compiler switches used for building shared libraries.

`--ldflags`

Print the value of the `LDFLAGS` variable that was used for building Postgres Pro. This shows linker switches.

`--ldflags_ex`

Print the value of the `LDFLAGS_EX` variable that was used for building Postgres Pro. This shows linker switches used for building executables only.

`--ldflags_sl`

Print the value of the `LDFLAGS_SL` variable that was used for building Postgres Pro. This shows linker switches used for building shared libraries only.

`--libs`

Print the value of the `LIBS` variable that was used for building Postgres Pro. This normally contains `-l` switches for external libraries linked into Postgres Pro.

`--version`

Print the PostgreSQL version on which Postgres Pro is based.

`--pgpro-version`

Print the version of Postgres Pro.

`--pgpro-edition`

Print the edition of Postgres Pro.

-?
--help

Show help about pg_config command line arguments, and exit.

If more than one option is given, the information is printed in that order, one item per line. If no options are given, all available information is printed, with labels.

Notes

The options `--docdir`, `--pkgincludedir`, `--localedir`, `--mandir`, `--sharedir`, `--sysconfdir`, `--cc`, `--cppflags`, `--cflags`, `--cflags_sl`, `--ldflags`, `--ldflags_sl`, and `--libs` were added in PostgreSQL 8.1. The option `--htmldir` was added in PostgreSQL 8.4. The option `--ldflags_ex` was added in PostgreSQL 9.0.

Example

To reproduce the build configuration of the current Postgres Pro installation, run the following command:

```
eval `./configure `pg_config --configure`
```

The output of `pg_config --configure` contains shell quotation marks so arguments with spaces are represented correctly. Therefore, using `eval` is required for proper results.

pg_dump

`pg_dump` — extract a Postgres Pro database into a script file or other archive file

Synopsis

```
pg_dump [connection-option...] [option...] [dbname]
```

Description

`pg_dump` is a utility for backing up a Postgres Pro database. It makes consistent backups even if the database is being used concurrently. `pg_dump` does not block other users accessing the database (readers or writers).

Note

Note `pg_dump` does not work with the distributed systems.

`pg_dump` only dumps a single database. To back up an entire cluster, or to back up global objects that are common to all databases in a cluster (such as roles and tablespaces), use [pg_dumpall](#).

Dumps can be output in script or archive file formats. Script dumps are plain-text files containing the SQL commands required to reconstruct the database to the state it was in at the time it was saved. To restore from such a script, feed it to [psql](#). Script files can be used to reconstruct the database even on other machines and other architectures; with some modifications, even on other SQL database products.

The alternative archive file formats must be used with [pg_restore](#) to rebuild the database. They allow `pg_restore` to be selective about what is restored, or even to reorder the items prior to being restored. The archive file formats are designed to be portable across architectures.

When used with one of the archive file formats and combined with `pg_restore`, `pg_dump` provides a flexible archival and transfer mechanism. `pg_dump` can be used to backup an entire database, then `pg_restore` can be used to examine the archive and/or select which parts of the database are to be restored. The most flexible output file formats are the “custom” format (`-Fc`) and the “directory” format (`-Fd`). They allow for selection and reordering of all archived items, support parallel restoration, and are compressed by default. The “directory” format is the only format that supports parallel dumps.

While running `pg_dump`, one should examine the output for any warnings (printed on standard error), especially in light of the limitations listed below.

Options

The following command-line options control the content and format of the output.

dbname

Specifies the name of the database to be dumped. If this is not specified, the environment variable `PGDATABASE` is used. If that is not set, the user name specified for the connection is used.

`-a`

`--data-only`

Dump only the data, not the schema (data definitions). Table data, large objects, and sequence values are dumped.

This option is similar to, but for historical reasons not identical to, specifying `--section=data`.

`-b`
`--large-objects`
`--blobs` (deprecated)

Include large objects in the dump. This is the default behavior except when `--schema`, `--table`, or `--schema-only` is specified. The `-b` switch is therefore only useful to add large objects to dumps where a specific schema or table has been requested. Note that large objects are considered data and therefore will be included when `--data-only` is used, but not when `--schema-only` is.

`-B`
`--no-large-objects`
`--no-blobs` (deprecated)

Exclude large objects in the dump.

When both `-b` and `-B` are given, the behavior is to output large objects, when data is being dumped, see the `-b` documentation.

`-c`
`--clean`

Output commands to `DROP` all the dumped database objects prior to outputting the commands for creating them. This option is useful when the restore is to overwrite an existing database. If any of the objects do not exist in the destination database, ignorable error messages will be reported during restore, unless `--if-exists` is also specified.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`-C`
`--create`

Begin the output with a command to create the database itself and reconnect to the created database. (With a script of this form, it doesn't matter which database in the destination installation you connect to before running the script.) If `--clean` is also specified, the script drops and recreates the target database before reconnecting to it.

With `--create`, the output also includes the database's comment if any, and any configuration variable settings that are specific to this database, that is, any `ALTER DATABASE ... SET ...` and `ALTER ROLE ... IN DATABASE ... SET ...` commands that mention this database. Access privileges for the database itself are also dumped, unless `--no-acl` is specified.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`-e pattern`
`--extension=pattern`

Dump only extensions matching *pattern*. When this option is not specified, all non-system extensions in the target database will be dumped. Multiple extensions can be selected by writing multiple `-e` switches. The *pattern* parameter is interpreted as a pattern according to the same rules used by `psql`'s `\d` commands (see [Patterns](#)), so multiple extensions can also be selected by writing wildcard characters in the pattern. When using wildcards, be careful to quote the pattern if needed to prevent the shell from expanding the wildcards.

Any configuration relation registered by `pg_extension_config_dump` is included in the dump if its extension is specified by `--extension`.

Note

When `-e` is specified, `pg_dump` makes no attempt to dump any other database objects that the selected extension(s) might depend upon. Therefore, there is no guarantee that the results of a specific-extension dump can be successfully restored by themselves into a clean database.

`-E encoding`
`--encoding=encoding`

Create the dump in the specified character set encoding. By default, the dump is created in the database encoding. (Another way to get the same result is to set the `PGCLIENTENCODING` environment variable to the desired dump encoding.) The supported encodings are described in [Section 23.3.1](#).

`-f file`
`--file=file`

Send output to the specified file. This parameter can be omitted for file based output formats, in which case the standard output is used. It must be given for the directory output format however, where it specifies the target directory instead of a file. In this case the directory is created by `pg_dump` and must not exist before.

`-F format`
`--format=format`

Selects the format of the output. *format* can be one of the following:

`p`
`plain`

Output a plain-text SQL script file (the default).

`c`
`custom`

Output a custom-format archive suitable for input into `pg_restore`. Together with the directory output format, this is the most flexible output format in that it allows manual selection and re-ordering of archived items during restore. This format is also compressed by default.

`d`
`directory`

Output a directory-format archive suitable for input into `pg_restore`. This will create a directory with one file for each table and large object being dumped, plus a so-called Table of Contents file describing the dumped objects in a machine-readable format that `pg_restore` can read. A directory format archive can be manipulated with standard Unix tools; for example, files in an uncompressed archive can be compressed with the `gzip`, `lz4`, or `zstd` tools. This format is compressed by default using `gzip` and also supports parallel dumps.

`t`
`tar`

Output a `tar`-format archive suitable for input into `pg_restore`. The `tar` format is compatible with the directory format: extracting a `tar`-format archive produces a valid directory-format archive. However, the `tar` format does not support compression. Also, when using `tar` format the relative order of table data items cannot be changed during restore.

`-j njobs`
`--jobs=njobs`

Run the dump in parallel by dumping *njobs* tables simultaneously. This option may reduce the time needed to perform the dump but it also increases the load on the database server. You can only use this option with the directory output format because this is the only output format where multiple processes can write their data at the same time.

`pg_dump` will open `njobs + 1` connections to the database, so make sure your `max_connections` setting is high enough to accommodate all connections.

Requesting exclusive locks on database objects while running a parallel dump could cause the dump to fail. The reason is that the `pg_dump` leader process requests shared locks ([ACCESS SHARE](#)) on the objects that the worker processes are going to dump later in order to make sure that nobody deletes them and makes them go away while the dump is running. If another client then requests an exclusive lock on a table, that lock will not be granted but will be queued waiting for the shared lock of the leader process to be released. Consequently any other access to the table will not be granted either and will queue after the exclusive lock request. This includes the worker process trying to dump the table. Without any precautions this would be a classic deadlock situation. To detect this conflict, the `pg_dump` worker process requests another shared lock using the `NOWAIT` option. If the worker process is not granted this shared lock, somebody else must have requested an exclusive lock in the meantime and there is no way to continue with the dump, so `pg_dump` has no choice but to abort the dump.

To perform a parallel dump, the database server needs to support synchronized snapshots, a feature that was introduced in PostgreSQL 9.2 for primary servers and 10 for standbys. With this feature, database clients can ensure they see the same data set even though they use different connections. `pg_dump -j` uses multiple database connections; it connects to the database once with the leader process and once again for each worker job. Without the synchronized snapshot feature, the different worker jobs wouldn't be guaranteed to see the same data in each connection, which could lead to an inconsistent backup.

```
-n pattern
--schema=pattern
```

Dump only schemas matching *pattern*; this selects both the schema itself, and all its contained objects. When this option is not specified, all non-system schemas in the target database will be dumped. Multiple schemas can be selected by writing multiple `-n` switches. The *pattern* parameter is interpreted as a pattern according to the same rules used by `psql's \d` commands (see [Patterns](#)), so multiple schemas can also be selected by writing wildcard characters in the pattern. When using wildcards, be careful to quote the pattern if needed to prevent the shell from expanding the wildcards; see [Examples](#) below.

Note

When `-n` is specified, `pg_dump` makes no attempt to dump any other database objects that the selected schema(s) might depend upon. Therefore, there is no guarantee that the results of a specific-schema dump can be successfully restored by themselves into a clean database.

Note

Non-schema objects such as large objects are not dumped when `-n` is specified. You can add large objects back to the dump with the `--large-objects` switch.

```
-N pattern
--exclude-schema=pattern
```

Do not dump any schemas matching *pattern*. The pattern is interpreted according to the same rules as for `-n`. `-N` can be given more than once to exclude schemas matching any of several patterns.

When both `-n` and `-N` are given, the behavior is to dump just the schemas that match at least one `-n` switch but no `-N` switches. If `-N` appears without `-n`, then schemas matching `-N` are excluded from what is otherwise a normal dump.

`-O`
`--no-owner`

Do not output commands to set ownership of objects to match the original database. By default, `pg_dump` issues `ALTER OWNER` or `SET SESSION AUTHORIZATION` statements to set ownership of created database objects. These statements will fail when the script is run unless it is started by a superuser (or the same user that owns all of the objects in the script). To make a script that can be restored by any user, but will give that user ownership of all the objects, specify `-O`.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`-R`
`--no-reconnect`

This option is obsolete but still accepted for backwards compatibility.

`-s`
`--schema-only`

Dump only the object definitions (schema), not data.

This option is the inverse of `--data-only`. It is similar to, but for historical reasons not identical to, specifying `--section=pre-data --section=post-data`.

(Do not confuse this with the `--schema` option, which uses the word “schema” in a different meaning.)

To exclude table data for only a subset of tables in the database, see `--exclude-table-data`.

`-S username`
`--superuser=username`

Specify the superuser user name to use when disabling triggers. This is relevant only if `--disable-triggers` is used. (Usually, it's better to leave this out, and instead start the resulting script as superuser.)

`-t pattern`
`--table=pattern`

Dump only tables with names matching *pattern*. Multiple tables can be selected by writing multiple `-t` switches. The *pattern* parameter is interpreted as a pattern according to the same rules used by `psql`'s `\d` commands (see [Patterns](#)), so multiple tables can also be selected by writing wildcard characters in the pattern. When using wildcards, be careful to quote the pattern if needed to prevent the shell from expanding the wildcards; see [Examples](#) below.

As well as tables, this option can be used to dump the definition of matching views, materialized views, foreign tables, and sequences. It will not dump the contents of views or materialized views, and the contents of foreign tables will only be dumped if the corresponding foreign server is specified with `--include-foreign-data`.

The `-n` and `-N` switches have no effect when `-t` is used, because tables selected by `-t` will be dumped regardless of those switches, and non-table objects will not be dumped.

Note

When `-t` is specified, `pg_dump` makes no attempt to dump any other database objects that the selected table(s) might depend upon. Therefore, there is no guarantee that the results of a specific-table dump can be successfully restored by themselves into a clean database.

`-T pattern`
`--exclude-table=pattern`

Do not dump any tables matching *pattern*. The pattern is interpreted according to the same rules as for `-t`. `-T` can be given more than once to exclude tables matching any of several patterns.

When both `-t` and `-T` are given, the behavior is to dump just the tables that match at least one `-t` switch but no `-T` switches. If `-T` appears without `-t`, then tables matching `-T` are excluded from what is otherwise a normal dump.

`-v`
`--verbose`

Specifies verbose mode. This will cause `pg_dump` to output detailed object comments and start/stop times to the dump file, and progress messages to standard error. Repeating the option causes additional debug-level messages to appear on standard error.

`-V`
`--version`

Print the `pg_dump` version and exit.

`-x`
`--no-privileges`
`--no-acl`

Prevent dumping of access privileges (grant/revoke commands).

`-Z level`
`-Z method[:detail]`
`--compress=level`
`--compress=method[:detail]`

Specify the compression method and/or the compression level to use. The compression method can be set to `gzip`, `lz4`, `zstd`, or `none` for no compression. A compression detail string can optionally be specified. If the detail string is an integer, it specifies the compression level. Otherwise, it should be a comma-separated list of items, each of the form `keyword` or `keyword=value`. Currently, the supported keywords are `level` and `long`.

If no compression level is specified, the default compression level will be used. If only a level is specified without mentioning an algorithm, `gzip` compression will be used if the level is greater than 0, and no compression will be used if the level is 0.

For the custom and directory archive formats, this specifies compression of individual table-data segments, and the default is to compress using `gzip` at a moderate level. For plain text output, setting a nonzero compression level causes the entire output file to be compressed, as though it had been fed through `gzip`, `lz4`, or `zstd`; but the default is not to compress. With `zstd` compression, `long` mode may improve the compression ratio, at the cost of increased memory use.

The tar archive format currently does not support compression at all.

`--binary-upgrade`

This option is for use by in-place upgrade utilities. Its use for other purposes is not recommended or supported. The behavior of the option may change in future releases without notice.

`--column-inserts`
`--attribute-inserts`

Dump data as `INSERT` commands with explicit column names (`INSERT INTO table (column, ...) VALUES ...`). This will make restoration very slow; it is mainly useful for making dumps that can

be loaded into non-Postgres Pro databases. Any error during restoring will cause only rows that are part of the problematic `INSERT` to be lost, rather than the entire table contents.

`--disable-dollar-quoting`

This option disables the use of dollar quoting for function bodies, and forces them to be quoted using SQL standard string syntax.

`--disable-triggers`

This option is relevant only when creating a data-only dump. It instructs `pg_dump` to include commands to temporarily disable triggers on the target tables while the data is restored. Use this if you have referential integrity checks or other triggers on the tables that you do not want to invoke during data restore.

Presently, the commands emitted for `--disable-triggers` must be done as superuser. So, you should also specify a superuser name with `-S`, or preferably be careful to start the resulting script as a superuser.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`--enable-large-mem-buffers`

Dump large `bytea` values exceeding 0.5 GB (up to 1 GB) and records with several `text` values exceeding 1 GB in total (up to 2 GB). This option cannot be used with `--inserts`, `--column-inserts`, or `--rows-per-insert`.

`--enable-row-security`

This option is relevant only when dumping the contents of a table which has row security. By default, `pg_dump` will set `row_security` to off, to ensure that all data is dumped from the table. If the user does not have sufficient privileges to bypass row security, then an error is thrown. This parameter instructs `pg_dump` to set `row_security` to on instead, allowing the user to dump the parts of the contents of the table that they have access to.

Note that if you use this option currently, you probably also want the dump be in `INSERT` format, as the `COPY FROM` during restore does not support row security.

`--exclude-extension=pattern`

Do not dump any extensions matching *pattern*. The pattern is interpreted according to the same rules as for `-e`. `--exclude-extension` can be given more than once to exclude extensions matching any of several patterns.

When both `-e` and `--exclude-extension` are given, the behavior is to dump just the extensions that match at least one `-e` switch but no `--exclude-extension` switches. If `--exclude-extension` appears without `-e`, then extensions matching `--exclude-extension` are excluded from what is otherwise a normal dump.

`--exclude-table-and-children=pattern`

This is the same as the `-T/--exclude-table` option, except that it also excludes any partitions or inheritance child tables of the table(s) matching the *pattern*.

`--exclude-table-data=pattern`

Do not dump data for any tables matching *pattern*. The pattern is interpreted according to the same rules as for `-t`. `--exclude-table-data` can be given more than once to exclude tables matching any of several patterns. This option is useful when you need the definition of a particular table even though you do not need the data in it.

To exclude data for all tables in the database, see `--schema-only`.

`--exclude-table-data-and-children=pattern`

This is the same as the `--exclude-table-data` option, except that it also excludes data of any partitions or inheritance child tables of the table(s) matching the *pattern*.

`--extra-float-digits=ndigits`

Use the specified value of `extra_float_digits` when dumping floating-point data, instead of the maximum available precision. Routine dumps made for backup purposes should not use this option.

`--filter=filename`

Specify a filename from which to read patterns for objects to include or exclude from the dump. The patterns are interpreted according to the same rules as the corresponding options: `-t/--table`, `--table-and-children`, `-T/--exclude-table`, and `--exclude-table-and-children` for tables, `-n/--schema` and `-N/--exclude-schema` for schemas, `--include-foreign-data` for data on foreign servers, `--exclude-table-data` and `--exclude-table-data-and-children` for table data, and `-e/--extension` and `--exclude-extension` for extensions. To read from `STDIN`, use `-` as the filename. The `--filter` option can be specified in conjunction with the above listed options for including or excluding objects, and can also be specified more than once for multiple filter files.

The file lists one object pattern per row, with the following format:

```
{ include | exclude } { extension | foreign_data | table | table_and_children |
table_data | table_data_and_children | schema } PATTERN
```

The first keyword specifies whether the objects matched by the pattern are to be included or excluded. The second keyword specifies the type of object to be filtered using the pattern:

- **extension:** extensions. This works like the `-e/--extension` or `--exclude-extension` option.
- **foreign_data:** data on foreign servers. This works like the `--include-foreign-data` option. This keyword can only be used with the `include` keyword.
- **table:** tables. This works like the `-t/--table` or `-T/--exclude-table` option.
- **table_and_children:** tables including any partitions or inheritance child tables. This works like the `--table-and-children` or `--exclude-table-and-children` option.
- **table_data:** table data of any tables matching *pattern*. This works like the `--exclude-table-data` option. This keyword can only be used with the `exclude` keyword.
- **table_data_and_children:** table data of any tables matching *pattern* as well as any partitions or inheritance children of the table(s). This works like the `--exclude-table-data-and-children` option. This keyword can only be used with the `exclude` keyword.
- **schema:** schemas. This works like the `-n/--schema` or `-N/--exclude-schema` option.

Lines starting with `#` are considered comments and ignored. Comments can be placed after an object pattern row as well. Blank lines are also ignored. See [Patterns](#) for how to perform quoting in patterns.

Example files are listed below in the [Examples](#) section.

`--if-exists`

Use `DROP ... IF EXISTS` commands to drop objects in `--clean` mode. This suppresses “does not exist” errors that might otherwise be reported. This option is not valid unless `--clean` is also specified.

`--include-foreign-data=foreignserver`

Dump the data for any foreign table with a foreign server matching *foreignserver* pattern. Multiple foreign servers can be selected by writing multiple `--include-foreign-data` switches. Also, the *foreignserver* parameter is interpreted as a pattern according to the same rules used by `psql's \d` commands (see [Patterns](#)), so multiple foreign servers can also be selected by writing wildcard characters in the pattern. When using wildcards, be careful to quote the pattern if needed to prevent

the shell from expanding the wildcards; see [Examples](#) below. The only exception is that an empty pattern is disallowed.

Note

Using wildcards in `--include-foreign-data` may result in access to unexpected foreign servers. Also, to use this option securely, make sure that the named server must have a trusted owner.

Note

When `--include-foreign-data` is specified, `pg_dump` does not check that the foreign table is writable. Therefore, there is no guarantee that the results of a foreign table dump can be successfully restored.

`--inserts`

Dump data as `INSERT` commands (rather than `COPY`). This will make restoration very slow; it is mainly useful for making dumps that can be loaded into non-Postgres Pro databases. Any error during restoring will cause only rows that are part of the problematic `INSERT` to be lost, rather than the entire table contents. Note that the restore might fail altogether if you have rearranged column order. The `--column-inserts` option is safe against column order changes, though even slower.

`--load-via-partition-root`

When dumping data for a table partition, make the `COPY` or `INSERT` statements target the root of the partitioning hierarchy that contains it, rather than the partition itself. This causes the appropriate partition to be re-determined for each row when the data is loaded. This may be useful when restoring data on a server where rows do not always fall into the same partitions as they did on the original server. That could happen, for example, if the partitioning column is of type `text` and the two systems have different definitions of the collation used to sort the partitioning column.

`--lock-wait-timeout=timeout`

Do not wait forever to acquire shared table locks at the beginning of the dump. Instead fail if unable to lock a table within the specified *timeout*. The timeout may be specified in any of the formats accepted by `SET statement_timeout`. (Allowed formats vary depending on the server version you are dumping from, but an integer number of milliseconds is accepted by all versions.)

`--no-comments`

Do not dump comments.

`--no-publications`

Do not dump publications.

`--no-security-labels`

Do not dump security labels.

`--no-subscriptions`

Do not dump subscriptions.

`--no-sync`

By default, `pg_dump` will wait for all files to be written safely to disk. This option causes `pg_dump` to return without waiting, which is faster, but means that a subsequent operating system crash can

leave the dump corrupt. Generally, this option is useful for testing but should not be used when dumping data from production installation.

`--no-table-access-method`

Do not output commands to select table access methods. With this option, all objects will be created with whichever table access method is the default during restore.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`--no-tablespaces`

Do not output commands to select tablespaces. With this option, all objects will be created in whichever tablespace is the default during restore.

This option is ignored when emitting an archive (non-text) output file. For the archive formats, you can specify the option when you call `pg_restore`.

`--no-toast-compression`

Do not output commands to set TOAST compression methods. With this option, all columns will be restored with the default compression setting.

`--no-unlogged-table-data`

Do not dump the contents of unlogged tables and sequences. This option has no effect on whether or not the table and sequence definitions (schema) are dumped; it only suppresses dumping the table and sequence data. Data in unlogged tables and sequences is always excluded when dumping from a standby server.

`--on-conflict-do-nothing`

Add `ON CONFLICT DO NOTHING` to `INSERT` commands. This option is not valid unless `--inserts`, `--column-inserts` or `--rows-per-insert` is also specified.

`--privileges-only`

Dump only access privileges (`GRANT/REVOKE` commands) and output commands executed by the security officer for schemas and their objects.

`--quote-all-identifiers`

Force quoting of all identifiers. This option is recommended when dumping a database from a server whose PostgreSQL major version is different from `pg_dump`'s, or when the output is intended to be loaded into a server of a different major version. By default, `pg_dump` quotes only identifiers that are reserved words in its own major version. This sometimes results in compatibility issues when dealing with servers of other versions that may have slightly different sets of reserved words. Using `--quote-all-identifiers` prevents such issues, at the price of a harder-to-read dump script.

`--rows-per-insert=nrows`

Dump data as `INSERT` commands (rather than `COPY`). Controls the maximum number of rows per `INSERT` command. The value specified must be a number greater than zero. Any error during restoring will cause only rows that are part of the problematic `INSERT` to be lost, rather than the entire table contents.

`--section=sectionname`

Only dump the named section. The section name can be `pre-data`, `data`, or `post-data`. This option can be specified more than once to select multiple sections. The default is to dump all sections.

The data section contains actual table data, large-object contents, and sequence values. Post-data items include definitions of indexes, triggers, rules, and constraints other than validated check constraints. Pre-data items include all other data definition items.

`--serializable-deferrable`

Use a `serializable` transaction for the dump, to ensure that the snapshot used is consistent with later database states; but do this by waiting for a point in the transaction stream at which no anomalies can be present, so that there isn't a risk of the dump failing or causing other transactions to roll back with a `serialization_failure`. See [Chapter 13](#) for more information about transaction isolation and concurrency control.

This option is not beneficial for a dump which is intended only for disaster recovery. It could be useful for a dump used to load a copy of the database for reporting or other read-only load sharing while the original database continues to be updated. Without it the dump may reflect a state which is not consistent with any serial execution of the transactions eventually committed. For example, if batch processing techniques are used, a batch may show as closed in the dump without all of the items which are in the batch appearing.

This option will make no difference if there are no read-write transactions active when `pg_dump` is started. If read-write transactions are active, the start of the dump may be delayed for an indeterminate length of time. Once running, performance with or without the switch is the same.

`--snapshot=snapshotname`

Use the specified synchronized snapshot when making a dump of the database (see [Table 9.98](#) for more details).

This option is useful when needing to synchronize the dump with a logical replication slot (see [Chapter 50](#)) or with a concurrent session.

In the case of a parallel dump, the snapshot name defined by this option is used rather than taking a new snapshot.

`--strict-names`

Require that each extension (`-e/--extension`), schema (`-n/--schema`) and table (`-t/--table`) pattern match at least one extension/schema/table in the database to be dumped. This also applies to filters used with `--filter`. Note that if none of the extension/schema/table patterns find matches, `pg_dump` will generate an error even without `--strict-names`.

This option has no effect on `--exclude-extension`, `-N/--exclude-schema`, `-T/--exclude-table`, or `--exclude-table-data`. An exclude pattern failing to match any objects is not considered an error.

`--sync-method=method`

When set to `fsync`, which is the default, `pg_dump --format=directory` will recursively open and synchronize all files in the archive directory.

On Linux, `syncfs` may be used instead to ask the operating system to synchronize the whole file system that contains the archive directory. See [recovery_init_sync_method](#) for information about the caveats to be aware of when using `syncfs`.

This option has no effect when `--no-sync` is used or `--format` is not set to `directory`.

`--table-and-children=pattern`

This is the same as the `-t/--table` option, except that it also includes any partitions or inheritance child tables of the table(s) matching the *pattern*.

`--use-set-session-authorization`

Output SQL-standard `SET SESSION AUTHORIZATION` commands instead of `ALTER OWNER` commands to determine object ownership. This makes the dump more standards-compatible, but depending on the history of the objects in the dump, might not restore properly. Also, a dump using `SET SESSION AUTHORIZATION` will certainly require superuser privileges to restore correctly, whereas `ALTER OWNER` requires lesser privileges.

-?
--help

Show help about pg_dump command line arguments, and exit.

The following command-line options control the database connection parameters.

-d *dbname*
--dbname=*dbname*

Specifies the name of the database to connect to. This is equivalent to specifying *dbname* as the first non-option argument on the command line. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

-h *host*
--host=*host*

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket. The default is taken from the PGHOST environment variable, if set, else a Unix domain socket connection is attempted.

-p *port*
--port=*port*

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections. Defaults to the PGPORT environment variable, if set, or a compiled-in default.

-U *username*
--username=*username*

User name to connect as.

-w
--no-password

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a .pgpass file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

-W
--password

Force pg_dump to prompt for a password before connecting to a database.

This option is never essential, since pg_dump will automatically prompt for a password if the server demands password authentication. However, pg_dump will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing -W to avoid the extra connection attempt.

--role=*rolename*

Specifies a role name to be used to create the dump. This option causes pg_dump to issue a SET ROLE *rolename* command after connecting to the database. It is useful when the authenticated user (specified by -U) lacks privileges needed by pg_dump, but can switch to a role with the required rights. Some installations have a policy against logging in directly as a superuser, and use of this option allows dumps to be made without violating the policy.

Environment

PGDATABASE
PGHOST
PGOPTIONS
PGPORT
PGUSER

Default connection parameters.

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

`pg_dump` internally executes `SELECT` statements. If you have problems running `pg_dump`, make sure you are able to select information from the database using, for example, `psql`. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

The database activity of `pg_dump` is normally collected by the cumulative statistics system. If this is undesirable, you can set parameter `track_counts` to `false` via `PGOPTIONS` or the `ALTER USER` command.

Notes

If your database cluster has any local additions to the `template1` database, be careful to restore the output of `pg_dump` into a truly empty database; otherwise you are likely to get errors due to duplicate definitions of the added objects. To make an empty database without any local additions, copy from `template0` not `template1`, for example:

```
CREATE DATABASE foo WITH TEMPLATE template0;
```

When a data-only dump is chosen and the option `--disable-triggers` is used, `pg_dump` emits commands to disable triggers on user tables before inserting the data, and then commands to re-enable them after the data has been inserted. If the restore is stopped in the middle, the system catalogs might be left in the wrong state.

The dump file produced by `pg_dump` does not contain the statistics used by the optimizer to make query planning decisions. Therefore, it is wise to run `ANALYZE` after restoring from a dump file to ensure optimal performance; see [Section 24.1.3](#) and [Section 24.1.6](#) for more information.

Because `pg_dump` is used to transfer data to newer versions of Postgres Pro, the output of `pg_dump` can be expected to load into Postgres Pro server versions newer than `pg_dump`'s version. `pg_dump` can also dump from Postgres Pro servers older than its own version. (Currently, servers back to version 9.2 are supported.) However, `pg_dump` cannot dump from Postgres Pro servers newer than its own major version; it will refuse to even try, rather than risk making an invalid dump. Also, it is not guaranteed that `pg_dump`'s output can be loaded into a server of an older major version — not even if the dump was taken from a server of that version. Loading a dump file into an older server may require manual editing of the dump file to remove syntax not understood by the older server. Use of the `--quote-all-identifiers` option is recommended in cross-version cases, as it can prevent problems arising from varying reserved-word lists in different PostgreSQL versions.

When dumping logical replication subscriptions, `pg_dump` will generate `CREATE SUBSCRIPTION` commands that use the `connect = false` option, so that restoring the subscription does not make remote connections for creating a replication slot or for initial table copy. That way, the dump can be restored without requiring network access to the remote servers. It is then up to the user to reactivate the subscriptions in a suitable way. If the involved hosts have changed, the connection information might have to be changed. It might also be appropriate to truncate the target tables before initiating a new full table copy. If users intend to copy initial data during refresh they must create the slot with `two_phase = false`. After the initial sync, the `two_phase` option will be automatically enabled by the subscriber if the subscription had been originally created with `two_phase = true` option.

It is generally recommended to use the `-X (--no-psqlrc)` option when restoring a database from a plain-text `pg_dump` script to ensure a clean restore process and prevent potential conflicts with non-default `psql` configurations.

Examples

To dump a database called `mydb` into an SQL-script file:

```
$ pg_dump mydb > db.sql
```

To reload such a script into a (freshly created) database named `newdb`:

```
$ psql -X -d newdb -f db.sql
```

To dump a database into a custom-format archive file:

```
$ pg_dump -Fc mydb > db.dump
```

To dump a database into a directory-format archive:

```
$ pg_dump -Fd mydb -f dumpdir
```

To dump a database into a directory-format archive in parallel with 5 worker jobs:

```
$ pg_dump -Fd mydb -j 5 -f dumpdir
```

To reload an archive file into a (freshly created) database named `newdb`:

```
$ pg_restore -d newdb db.dump
```

To reload an archive file into the same database it was dumped from, discarding the current contents of that database:

```
$ pg_restore -d postgres --clean --create db.dump
```

To dump a single table named `mytab`:

```
$ pg_dump -t mytab mydb > db.sql
```

To dump all tables whose names start with `emp` in the `detroit` schema, except for the table named `employee_log`:

```
$ pg_dump -t 'detroit.emp*' -T detroit.employee_log mydb > db.sql
```

To dump all schemas whose names start with `east` or `west` and end in `gsm`, excluding any schemas whose names contain the word `test`:

```
$ pg_dump -n 'east*gsm' -n 'west*gsm' -N '*test*' mydb > db.sql
```

The same, using regular expression notation to consolidate the switches:

```
$ pg_dump -n '(east|west)*gsm' -N '*test*' mydb > db.sql
```

To dump all database objects except for tables whose names begin with `ts_`:

```
$ pg_dump -T 'ts_*' mydb > db.sql
```

To specify an upper-case or mixed-case name in `-t` and related switches, you need to double-quote the name; else it will be folded to lower case (see [Patterns](#)). But double quotes are special to the shell, so in turn they must be quoted. Thus, to dump a single table with a mixed-case name, you need something like

```
$ pg_dump -t "\"MixedCaseName\"" mydb > mytab.sql
```

To dump all tables whose names start with `mytable`, except for table `mytable2`, specify a filter file `filter.txt` like:

```
include table mytable*
exclude table mytable2
```

```
$ pg_dump --filter=filter.txt mydb > db.sql
```

See Also

[pg_dumpall](#), [pg_restore](#), [psql](#)

pg_dumpall

pg_dumpall — extract a Postgres Pro database cluster into a script file

Synopsis

```
pg_dumpall [connection-option...] [option...]
```

Description

pg_dumpall is a utility for writing out (“dumping”) all Postgres Pro databases of a cluster into one script file. The script file contains SQL commands that can be used as input to [psql](#) to restore the databases. It does this by calling [pg_dump](#) for each database in the cluster. pg_dumpall also dumps global objects that are common to all databases, namely database roles, tablespaces, and privilege grants for configuration parameters. (pg_dump does not save these objects.)

Note

Note pg_dumpall does not work with the distributed systems.

Since pg_dumpall reads tables from all databases you will most likely have to connect as a database superuser in order to produce a complete dump. Also you will need superuser privileges to execute the saved script in order to be allowed to add roles and create databases.

The SQL script will be written to the standard output. Use the `-f/--file` option or shell operators to redirect it into a file.

pg_dumpall needs to connect several times to the Postgres Pro server (once per database). If you use password authentication it will ask for a password each time. It is convenient to have a `~/.pgpass` file in such cases. See [Section 35.16](#) for more information.

Options

The following command-line options control the content and format of the output.

`-a`
`--data-only`

Dump only the data, not the schema (data definitions).

`-c`
`--clean`

Emit SQL commands to `DROP` all the dumped databases, roles, and tablespaces before recreating them. This option is useful when the restore is to overwrite an existing cluster. If any of the objects do not exist in the destination cluster, ignorable error messages will be reported during restore, unless `--if-exists` is also specified.

`-E encoding`
`--encoding=encoding`

Create the dump in the specified character set encoding. By default, the dump is created in the database encoding. (Another way to get the same result is to set the `PGCLIENTENCODING` environment variable to the desired dump encoding.)

`-f filename`
`--file=filename`

Send output to the specified file. If this is omitted, the standard output is used.

`--filter=filename`

Specify a filename from which to read patterns for databases excluded from the dump. The patterns are interpreted according to the same rules as `--exclude-database`. To read from `STDIN`, use `-` as the filename. The `--filter` option can be specified in conjunction with `--exclude-database` for excluding databases, and can also be specified more than once for multiple filter files.

The file lists one database pattern per row, with the following format:

```
exclude database PATTERN
```

Lines starting with `#` are considered comments and ignored. Comments can be placed after an object pattern row as well. Blank lines are also ignored. See [Patterns](#) for how to perform quoting in patterns.

`-g`

`--globals-only`

Dump only global objects (roles and tablespaces), no databases.

`-O`

`--no-owner`

Do not output commands to set ownership of objects to match the original database. By default, `pg_dumpall` issues `ALTER OWNER` or `SET SESSION AUTHORIZATION` statements to set ownership of created schema elements. These statements will fail when the script is run unless it is started by a superuser (or the same user that owns all of the objects in the script). To make a script that can be restored by any user, but will give that user ownership of all the objects, specify `-O`.

`-r`

`--roles-only`

Dump only roles, no databases or tablespaces.

`-s`

`--schema-only`

Dump only the object definitions (schema), not data.

`-S username`

`--superuser=username`

Specify the superuser user name to use when disabling triggers. This is relevant only if `--disable-triggers` is used. (Usually, it's better to leave this out, and instead start the resulting script as superuser.)

`-t`

`--tablespaces-only`

Dump only tablespaces, no databases or roles.

`-v`

`--verbose`

Specifies verbose mode. This will cause `pg_dumpall` to output start/stop times to the dump file, and progress messages to standard error. Repeating the option causes additional debug-level messages to appear on standard error. The option is also passed down to `pg_dump`.

`-V`

`--version`

Print the `pg_dumpall` version and exit.

`-x`

`--no-privileges`

`--no-acl`

Prevent dumping of access privileges (grant/revoke commands).

`--binary-upgrade`

This option is for use by in-place upgrade utilities. Its use for other purposes is not recommended or supported. The behavior of the option may change in future releases without notice.

`--column-inserts`

`--attribute-inserts`

Dump data as `INSERT` commands with explicit column names (`INSERT INTO table (column, ...) VALUES ...`). This will make restoration very slow; it is mainly useful for making dumps that can be loaded into non-Postgres Pro databases.

`--disable-dollar-quoting`

This option disables the use of dollar quoting for function bodies, and forces them to be quoted using SQL standard string syntax.

`--disable-triggers`

This option is relevant only when creating a data-only dump. It instructs `pg_dumpall` to include commands to temporarily disable triggers on the target tables while the data is restored. Use this if you have referential integrity checks or other triggers on the tables that you do not want to invoke during data restore.

Presently, the commands emitted for `--disable-triggers` must be done as superuser. So, you should also specify a superuser name with `-s`, or preferably be careful to start the resulting script as a superuser.

`--enable-large-mem-buffers`

Dump large bytea values exceeding 0.5 GB (up to 1 GB) and records with several text values exceeding 1 GB in total (up to 2 GB). This option cannot be used with `--inserts`, `--column-inserts`, or `--rows-per-insert`.

`--exclude-database=pattern`

Do not dump databases whose name matches *pattern*. Multiple patterns can be excluded by writing multiple `--exclude-database` switches. The *pattern* parameter is interpreted as a pattern according to the same rules used by `psql`'s `\d` commands (see [Patterns](#)), so multiple databases can also be excluded by writing wildcard characters in the pattern. When using wildcards, be careful to quote the pattern if needed to prevent shell wildcard expansion.

`--extra-float-digits=ndigits`

Use the specified value of `extra_float_digits` when dumping floating-point data, instead of the maximum available precision. Routine dumps made for backup purposes should not use this option.

`--if-exists`

Use `DROP ... IF EXISTS` commands to drop objects in `--clean` mode. This suppresses “does not exist” errors that might otherwise be reported. This option is not valid unless `--clean` is also specified.

`--inserts`

Dump data as `INSERT` commands (rather than `COPY`). This will make restoration very slow; it is mainly useful for making dumps that can be loaded into non-Postgres Pro databases. Note that the restore might fail altogether if you have rearranged column order. The `--column-inserts` option is safer, though even slower.

`--load-via-partition-root`

When dumping data for a table partition, make the `COPY` or `INSERT` statements target the root of the partitioning hierarchy that contains it, rather than the partition itself. This causes the appropriate

partition to be re-determined for each row when the data is loaded. This may be useful when restoring data on a server where rows do not always fall into the same partitions as they did on the original server. That could happen, for example, if the partitioning column is of type text and the two systems have different definitions of the collation used to sort the partitioning column.

`--lock-wait-timeout=timeout`

Do not wait forever to acquire shared table locks at the beginning of the dump. Instead, fail if unable to lock a table within the specified *timeout*. The timeout may be specified in any of the formats accepted by `SET statement_timeout`.

`--no-comments`

Do not dump comments.

`--no-publications`

Do not dump publications.

`--no-role-passwords`

Do not dump passwords for roles. When restored, roles will have a null password, and password authentication will always fail until the password is set. Since password values aren't needed when this option is specified, the role information is read from the catalog view `pg_roles` instead of `pg_authid`. Therefore, this option also helps if access to `pg_authid` is restricted by some security policy.

`--no-security-labels`

Do not dump security labels.

`--no-subscriptions`

Do not dump subscriptions.

`--no-sync`

By default, `pg_dumpall` will wait for all files to be written safely to disk. This option causes `pg_dumpall` to return without waiting, which is faster, but means that a subsequent operating system crash can leave the dump corrupt. Generally, this option is useful for testing but should not be used when dumping data from production installation.

`--no-table-access-method`

Do not output commands to select table access methods. With this option, all objects will be created with whichever table access method is the default during restore.

`--no-tablespaces`

Do not output commands to create tablespaces nor select tablespaces for objects. With this option, all objects will be created in whichever tablespace is the default during restore.

`--no-toast-compression`

Do not output commands to set TOAST compression methods. With this option, all columns will be restored with the default compression setting.

`--no-unlogged-table-data`

Do not dump the contents of unlogged tables. This option has no effect on whether or not the table definitions (schema) are dumped; it only suppresses dumping the table data.

`--on-conflict-do-nothing`

Add `ON CONFLICT DO NOTHING` to `INSERT` commands. This option is not valid unless `--inserts` or `--column-inserts` is also specified.

`--quote-all-identifiers`

Force quoting of all identifiers. This option is recommended when dumping a database from a server whose PostgreSQL major version is different from `pg_dumpall`'s, or when the output is intended to be loaded into a server of a different major version. By default, `pg_dumpall` quotes only identifiers that are reserved words in its own major version. This sometimes results in compatibility issues when dealing with servers of other versions that may have slightly different sets of reserved words. Using `--quote-all-identifiers` prevents such issues, at the price of a harder-to-read dump script.

`--rows-per-insert=nrows`

Dump data as `INSERT` commands (rather than `COPY`). Controls the maximum number of rows per `INSERT` command. The value specified must be a number greater than zero. Any error during restoring will cause only rows that are part of the problematic `INSERT` to be lost, rather than the entire table contents.

`--use-set-session-authorization`

Output SQL-standard `SET SESSION AUTHORIZATION` commands instead of `ALTER OWNER` commands to determine object ownership. This makes the dump more standards compatible, but depending on the history of the objects in the dump, might not restore properly.

`-?`

`--help`

Show help about `pg_dumpall` command line arguments, and exit.

The following command-line options control the database connection parameters.

`-d connstr`

`--dbname=connstr`

Specifies parameters used to connect to the server, as a [connection string](#); these will override any conflicting command line options.

The option is called `--dbname` for consistency with other client applications, but because `pg_dumpall` needs to connect to many databases, the database name in the connection string will be ignored. Use the `-l` option to specify the name of the database used for the initial connection, which will dump global objects and discover what other databases should be dumped.

`-h host`

`--host=host`

Specifies the host name of the machine on which the database server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket. The default is taken from the `PGHOST` environment variable, if set, else a Unix domain socket connection is attempted.

`-l dbname`

`--database=dbname`

Specifies the name of the database to connect to for dumping global objects and discovering what other databases should be dumped. If not specified, the `postgres` database will be used, and if that does not exist, `template1` will be used.

`-p port`

`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections. Defaults to the `PGPORT` environment variable, if set, or a compiled-in default.

`-U username`

`--username=username`

User name to connect as.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force `pg_dumpall` to prompt for a password before connecting to a database.

This option is never essential, since `pg_dumpall` will automatically prompt for a password if the server demands password authentication. However, `pg_dumpall` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

Note that the password prompt will occur again for each database to be dumped. Usually, it's better to set up a `~/.pgpass` file than to rely on manual password entry.

`--role=rolename`

Specifies a role name to be used to create the dump. This option causes `pg_dumpall` to issue a `SET ROLE rolename` command after connecting to the database. It is useful when the authenticated user (specified by `-U`) lacks privileges needed by `pg_dumpall`, but can switch to a role with the required rights. Some installations have a policy against logging in directly as a superuser, and use of this option allows dumps to be made without violating the policy.

Environment

`PGHOST`
`PGOPTIONS`
`PGPORT`
`PGUSER`

Default connection parameters

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by `libpq` (see [Section 35.15](#)).

Notes

Since `pg_dumpall` calls `pg_dump` internally, some diagnostic messages will refer to `pg_dump`.

The `--clean` option can be useful even when your intention is to restore the dump script into a fresh cluster. Use of `--clean` authorizes the script to drop and re-create the built-in `postgres` and `template1` databases, ensuring that those databases will retain the same properties (for instance, locale and encoding) that they had in the source cluster. Without the option, those databases will retain their existing database-level properties, as well as any pre-existing contents.

Once restored, it is wise to run `ANALYZE` on each database so the optimizer has useful statistics. You can also run `vacuumdb -a -z` to analyze all databases.

The dump script should not be expected to run completely without errors. In particular, because the script will issue `CREATE ROLE` for every role existing in the source cluster, it is certain to get a “role already exists” error for the bootstrap superuser, unless the destination cluster was initialized with a different bootstrap superuser name. This error is harmless and should be ignored. Use of the `--clean`

option is likely to produce additional harmless error messages about non-existent objects, although you can minimize those by adding `--if-exists`.

`pg_dumpall` requires all needed tablespace directories to exist before the restore; otherwise, database creation will fail for databases in non-default locations.

It is generally recommended to use the `-X` (`--no-psqlrc`) option when restoring a database from a `pg_dumpall` script to ensure a clean restore process and prevent potential conflicts with non-default `psql` configurations. Additionally, because the `pg_dumpall` script may include `psql` meta-commands, it may be incompatible with clients other than `psql`.

Examples

To dump all databases:

```
$ pg_dumpall > db.out
```

To restore database(s) from this file, you can use:

```
$ psql -X -f db.out -d postgres
```

It is not important which database you connect to here since the script file created by `pg_dumpall` will contain the appropriate commands to create and connect to the saved databases. An exception is that if you specified `--clean`, you must connect to the `postgres` database initially; the script will attempt to drop other databases immediately, and that will fail for the database you are connected to.

See Also

Check [pg_dump](#) for details on possible error conditions.

pg_isready

pg_isready — check the connection status of a Postgres Pro server

Synopsis

```
pg_isready [connection-option...] [option...]
```

Description

pg_isready is a utility for checking the connection status of a Postgres Pro database server. The exit status specifies the result of the connection check.

Options

```
-d dbname  
--dbname=dbname
```

Specifies the name of the database to connect to. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

```
-h hostname  
--host=hostname
```

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix-domain socket.

```
-p port  
--port=port
```

Specifies the TCP port or the local Unix-domain socket file extension on which the server is listening for connections. Defaults to the value of the `PGPORT` environment variable or, if not set, to the port specified at compile time, usually 5432.

```
-q  
--quiet
```

Do not display status message. This is useful when scripting.

```
-t seconds  
--timeout=seconds
```

The maximum number of seconds to wait when attempting connection before returning that the server is not responding. Setting to 0 disables. The default is 3 seconds.

```
-U username  
--username=username
```

Connect to the database as the user *username* instead of the default.

```
-V  
--version
```

Print the pg_isready version and exit.

```
-?  
--help
```

Show help about pg_isready command line arguments, and exit.

Exit Status

`pg_isready` returns 0 to the shell if the server is accepting connections normally, 1 if the server is rejecting connections (for example during startup), 2 if there was no response to the connection attempt, and 3 if no attempt was made (for example due to invalid parameters).

Environment

`pg_isready`, like most other Postgres Pro utilities, also uses the environment variables supported by `libpq` (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

It is not necessary to supply correct user name, password, or database name values to obtain the server status; however, if incorrect values are provided, the server will log a failed connection attempt.

Examples

Standard Usage:

```
$ pg_isready
/tmp:5432 - accepting connections
$ echo $?
0
```

Running with connection parameters to a Postgres Pro cluster in startup:

```
$ pg_isready -h localhost -p 5433
localhost:5433 - rejecting connections
$ echo $?
1
```

Running with connection parameters to a non-responsive Postgres Pro cluster:

```
$ pg_isready -h someremotehost
someremotehost:5432 - no response
$ echo $?
2
```

pg_receivewal

pg_receivewal — stream write-ahead logs from a Postgres Pro server

Synopsis

```
pg_receivewal [option...]
```

Description

pg_receivewal is used to stream the write-ahead log from a running Postgres Pro cluster. The write-ahead log is streamed using the streaming replication protocol, and is written to a local directory of files. This directory can be used as the archive location for doing a restore using point-in-time recovery (see [Section 25.1](#)).

pg_receivewal streams the write-ahead log in real time as it's being generated on the server, and does not wait for segments to complete like [archive command](#) and [archive library](#) do. For this reason, it is not necessary to set [archive_timeout](#) when using pg_receivewal.

Unlike the WAL receiver of a Postgres Pro standby server, pg_receivewal by default flushes WAL data only when a WAL file is closed. The option `--synchronous` must be specified to flush WAL data in real time. Since pg_receivewal does not apply WAL, you should not allow it to become a synchronous standby when [synchronous_commit](#) equals `remote_apply`. If it does, it will appear to be a standby that never catches up, and will cause transaction commits to block. To avoid this, you should either configure an appropriate value for [synchronous_standby_names](#), or specify `application_name` for pg_receivewal that does not match it, or change the value of `synchronous_commit` to something other than `remote_apply`.

The write-ahead log is streamed over a regular Postgres Pro connection and uses the replication protocol. The connection must be made with a user having `REPLICATION` permissions (see [Section 21.2](#)) or a superuser, and `pg_hba.conf` must permit the replication connection. The server must also be configured with [max_wal_senders](#) set high enough to leave at least one session available for the stream.

The starting point of the write-ahead log streaming is calculated when pg_receivewal starts:

1. First, scan the directory where the WAL segment files are written and find the newest completed segment file, using as the starting point the beginning of the next WAL segment file.
2. If a starting point cannot be calculated with the previous method, and if a replication slot is used, an extra `READ_REPLICATION_SLOT` command is issued to retrieve the slot's `restart_lsn` to use as the starting point. This option is only available when streaming write-ahead logs from Postgres Pro 15 and up.
3. If a starting point cannot be calculated with the previous method, the latest WAL flush location is used as reported by the server from an `IDENTIFY_SYSTEM` command.

If the connection is lost, or if it cannot be initially established, with a non-fatal error, pg_receivewal will retry the connection indefinitely, and reestablish streaming as soon as possible. To avoid this behavior, use the `-n` parameter.

In the absence of fatal errors, pg_receivewal will run until terminated by the SIGINT (**Control+C**) or SIGTERM signal.

Options

```
-D directory  
--directory=directory
```

Directory to write the output to.

This parameter is required.

`-E lsn`

`--endpos=lsn`

Automatically stop replication and exit with normal exit status 0 when receiving reaches the specified LSN.

If there is a record with LSN exactly equal to *lsn*, the record will be processed.

`--if-not-exists`

Do not error out when `--create-slot` is specified and a slot with the specified name already exists.

`-n`

`--no-loop`

Don't loop on connection errors. Instead, exit right away with an error.

`--no-sync`

This option causes `pg_receivewal` to not force WAL data to be flushed to disk. This is faster, but means that a subsequent operating system crash can leave the WAL segments corrupt. Generally, this option is useful for testing but should not be used when doing WAL archiving on a production deployment.

This option is incompatible with `--synchronous`.

`-s interval`

`--status-interval=interval`

Specifies the number of seconds between status packets sent back to the server. This allows for easier monitoring of the progress from server. A value of zero disables the periodic status updates completely, although an update will still be sent when requested by the server, to avoid timeout disconnect. The default value is 10 seconds.

`-S slotname`

`--slot=slotname`

Require `pg_receivewal` to use an existing replication slot (see [Section 26.2.6](#)). When this option is used, `pg_receivewal` will report a flush position to the server, indicating when each segment has been synchronized to disk so that the server can remove that segment if it is not otherwise needed.

When the replication client of `pg_receivewal` is configured on the server as a synchronous standby, then using a replication slot will report the flush position to the server, but only when a WAL file is closed. Therefore, that configuration will cause transactions on the primary to wait for a long time and effectively not work satisfactorily. The option `--synchronous` (see below) must be specified in addition to make this work correctly.

`--synchronous`

Flush the WAL data to disk immediately after it has been received. Also send a status packet back to the server immediately after flushing, regardless of `--status-interval`.

This option should be specified if the replication client of `pg_receivewal` is configured on the server as a synchronous standby, to ensure that timely feedback is sent to the server.

`-v`

`--verbose`

Enables verbose mode.

`-Z level`

`-Z method[:detail]`

`--compress=level`

`--compress=method[:detail]`

Enables compression of write-ahead logs.

The compression method can be set to `gzip`, `lz4` (if Postgres Pro was compiled with `--with-lz4`) or `none` for no compression. A compression detail string can optionally be specified. If the detail string is an integer, it specifies the compression level. Otherwise, it should be a comma-separated list of items, each of the form *keyword* or *keyword=value*. Currently, the only supported keyword is `level`.

If no compression level is specified, the default compression level will be used. If only a level is specified without mentioning an algorithm, `gzip` compression will be used if the level is greater than 0, and no compression will be used if the level is 0.

The suffix `.gz` will automatically be added to all filenames when using `gzip`, and the suffix `.lz4` is added when using `lz4`.

The following command-line options control the database connection parameters.

```
-d connstr
--dbname=connstr
```

Specifies parameters used to connect to the server, as a [connection string](#); these will override any conflicting command line options.

This option is called `--dbname` for consistency with other client applications, but because `pg_receivewal` doesn't connect to any particular database in the cluster, any database name included in the connection string will be ignored by the server. However, a database name supplied that way overrides the default database name (`replication`) for purposes of looking up the replication connection's password in `~/.pgpass`. Similarly, middleware or proxies used in connecting to Postgres Pro might utilize the name for purposes such as connection routing.

```
-h host
--host=host
```

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket. The default is taken from the `PGHOST` environment variable, if set, else a Unix domain socket connection is attempted.

```
-p port
--port=port
```

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections. Defaults to the `PGPORT` environment variable, if set, or a compiled-in default.

```
-U username
--username=username
```

User name to connect as.

```
-w
--no-password
```

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

```
-W
--password
```

Force `pg_receivewal` to prompt for a password before connecting to a database.

This option is never essential, since `pg_receivewal` will automatically prompt for a password if the server demands password authentication. However, `pg_receivewal` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`pg_receivewal` can perform one of the two following actions in order to control physical replication slots:

`--create-slot`

Create a new physical replication slot with the name specified in `--slot`, then exit.

`--drop-slot`

Drop the replication slot with the name specified in `--slot`, then exit.

Other options are also available:

`-V`

`--version`

Print the `pg_receivewal` version and exit.

`-?`

`--help`

Show help about `pg_receivewal` command line arguments, and exit.

Exit Status

`pg_receivewal` will exit with status 0 when terminated by the SIGINT or SIGTERM signal. (That is the normal way to end it. Hence it is not an error.) For fatal errors or other signals, the exit status will be nonzero.

Environment

This utility, like most other Postgres Pro utilities, uses the environment variables supported by libpq (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

When using `pg_receivewal` instead of [archive_command](#) or [archive_library](#) as the main WAL backup method, it is strongly recommended to use replication slots. Otherwise, the server is free to recycle or remove write-ahead log files before they are backed up, because it does not have any information, either from [archive_command](#) or [archive_library](#) or the replication slots, about how far the WAL stream has been archived. Note, however, that a replication slot will fill up the server's disk space if the receiver does not keep up with fetching the WAL data.

`pg_receivewal` will preserve group permissions on the received WAL files if group permissions are enabled on the source cluster.

Examples

To stream the write-ahead log from the server at `mydbserver` and store it in the local directory `/usr/local/pgsql/archive`:

```
$ pg_receivewal -h mydbserver -D /usr/local/pgsql/archive
```

pg_recvlogical

pg_recvlogical — control Postgres Pro logical decoding streams

Synopsis

```
pg_recvlogical [option...]
```

Description

pg_recvlogical controls logical decoding replication slots and streams data from such replication slots.

It creates a replication-mode connection, so it is subject to the same constraints as [pg_receivewal](#), plus those for logical replication (see [Chapter 50](#)).

pg_recvlogical has no equivalent to the logical decoding SQL interface's peek and get modes. It sends replay confirmations for data lazily as it receives it and on clean exit. To examine pending data on a slot without consuming it, use [pg_logical_slot_peek_changes](#).

In the absence of fatal errors, pg_recvlogical will run until terminated by the SIGINT (**Control+C**) or SIGTERM signal.

Options

At least one of the following options must be specified to select an action:

`--create-slot`

Create a new logical replication slot with the name specified by `--slot`, using the output plugin specified by `--plugin`, for the database specified by `--dbname`.

The `--two-phase` can be specified with `--create-slot` to enable decoding of prepared transactions.

`--drop-slot`

Drop the replication slot with the name specified by `--slot`, then exit.

`--start`

Begin streaming changes from the logical replication slot specified by `--slot`, continuing until terminated by a signal. If the server side change stream ends with a server shutdown or disconnect, retry in a loop unless `--no-loop` is specified.

The stream format is determined by the output plugin specified when the slot was created.

The connection must be to the same database used to create the slot.

`--create-slot` and `--start` can be specified together. `--drop-slot` cannot be combined with another action.

The following command-line options control the location and format of the output and other replication behavior:

`-E lsn`

`--endpos=lsn`

In `--start` mode, automatically stop replication and exit with normal exit status 0 when receiving reaches the specified LSN. If specified when not in `--start` mode, an error is raised.

If there's a record with LSN exactly equal to `lsn`, the record will be output.

The `--endpos` option is not aware of transaction boundaries and may truncate output partway through a transaction. Any partially output transaction will not be consumed and will be replayed again when the slot is next read from. Individual messages are never truncated.

`-f filename`
`--file=filename`

Write received and decoded transaction data into this file. Use `-` for stdout.

`-F interval_seconds`
`--fsync-interval=interval_seconds`

Specifies how often `pg_recvlogical` should issue `fsync()` calls to ensure the output file is safely flushed to disk.

The server will occasionally request the client to perform a flush and report the flush position to the server. This setting is in addition to that, to perform flushes more frequently.

Specifying an interval of 0 disables issuing `fsync()` calls altogether, while still reporting progress to the server. In this case, data could be lost in the event of a crash.

`-I lsn`
`--startpos=lsn`

In `--start` mode, start replication from the given LSN. For details on the effect of this, see the documentation in [Chapter 50](#) and [Section 56.4](#). Ignored in other modes.

`--if-not-exists`

Do not error out when `--create-slot` is specified and a slot with the specified name already exists.

`-n`
`--no-loop`

When the connection to the server is lost, do not retry in a loop, just exit.

`-o name[=value]`
`--option=name[=value]`

Pass the option *name* to the output plugin with, if specified, the option value *value*. Which options exist and their effects depends on the used output plugin.

`-P plugin`
`--plugin=plugin`

When creating a slot, use the specified logical decoding output plugin. See [Chapter 50](#). This option has no effect if the slot already exists.

`-s interval_seconds`
`--status-interval=interval_seconds`

This option has the same effect as the option of the same name in [pg_receivewal](#). See the description there.

`-S slot_name`
`--slot=slot_name`

In `--start` mode, use the existing logical replication slot named *slot_name*. In `--create-slot` mode, create the slot with this name. In `--drop-slot` mode, delete the slot with this name.

`-t`
`--two-phase`

Enables decoding of prepared transactions. This option may only be specified with `--create-slot`.

`-v`
`--verbose`

Enables verbose mode.

The following command-line options control the database connection parameters.

`-d dbname`
`--dbname=dbname`

The database to connect to. See the description of the actions for what this means in detail. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options. Defaults to the user name.

`-h hostname-or-ip`
`--host=hostname-or-ip`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket. The default is taken from the `PGHOST` environment variable, if set, else a Unix domain socket connection is attempted.

`-p port`
`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections. Defaults to the `PGPORT` environment variable, if set, or a compiled-in default.

`-U user`
`--username=user`

User name to connect as. Defaults to current operating system user name.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force `pg_recvlogical` to prompt for a password before connecting to a database.

This option is never essential, since `pg_recvlogical` will automatically prompt for a password if the server demands password authentication. However, `pg_recvlogical` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

The following additional options are available:

`-V`
`--version`

Print the `pg_recvlogical` version and exit.

`-?`
`--help`

Show help about `pg_recvlogical` command line arguments, and exit.

Exit Status

`pg_recvlogical` will exit with status 0 when terminated by the `SIGINT` or `SIGTERM` signal. (That is the normal way to end it. Hence it is not an error.) For fatal errors or other signals, the exit status will be nonzero.

Environment

This utility, like most other Postgres Pro utilities, uses the environment variables supported by libpq (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

`pg_recvlogical` will preserve group permissions on the received WAL files if group permissions are enabled on the source cluster.

Examples

See [Section 50.1](#) for an example.

See Also

[pg_receivewal](#)

pg_restore

pg_restore — restore a Postgres Pro database from an archive file created by pg_dump

Synopsis

```
pg_restore [connection-option...] [option...] [filename]
```

Description

pg_restore is a utility for restoring a Postgres Pro database from an archive created by [pg_dump](#) in one of the non-plain-text formats. It will issue the commands necessary to reconstruct the database to the state it was in at the time it was saved. The archive files also allow pg_restore to be selective about what is restored, or even to reorder the items prior to being restored. The archive files are designed to be portable across architectures.

pg_restore can operate in two modes. If a database name is specified, pg_restore connects to that database and restores archive contents directly into the database. Otherwise, a script containing the SQL commands necessary to rebuild the database is created and written to a file or standard output. This script output is equivalent to the plain text output format of pg_dump. Some of the options controlling the output are therefore analogous to pg_dump options.

Obviously, pg_restore cannot restore information that is not present in the archive file. For instance, if the archive was made using the “dump data as INSERT commands” option, pg_restore will not be able to load the data using COPY statements.

Options

pg_restore accepts the following command line arguments.

filename

Specifies the location of the archive file (or directory, for a directory-format archive) to be restored. If not specified, the standard input is used.

-a

--data-only

Restore only the data, not the schema (data definitions). Table data, large objects, and sequence values are restored, if present in the archive.

This option is similar to, but for historical reasons not identical to, specifying --section=data.

-c

--clean

Before restoring database objects, issue commands to DROP all the objects that will be restored. This option is useful for overwriting an existing database. If any of the objects do not exist in the destination database, ignorable error messages will be reported, unless --if-exists is also specified.

-C

--create

Create the database before restoring into it. If --clean is also specified, drop and recreate the target database before connecting to it.

With --create, pg_restore also restores the database's comment if any, and any configuration variable settings that are specific to this database, that is, any ALTER DATABASE ... SET ... and ALTER

ROLE ... IN DATABASE ... SET ... commands that mention this database. Access privileges for the database itself are also restored, unless `--no-acl` is specified.

When this option is used, the database named with `-d` is used only to issue the initial `DROP DATABASE` and `CREATE DATABASE` commands. All data is restored into the database name that appears in the archive.

`-d dbname`
`--dbname=dbname`

Connect to database *dbname* and restore directly into the database. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

`-e`
`--exit-on-error`

Exit if an error is encountered while sending SQL commands to the database. The default is to continue and to display a count of errors at the end of the restoration.

`-f filename`
`--file=filename`

Specify output file for generated script, or for the listing when used with `-l`. Use `-` for stdout.

`--filter=filename`

Specify a filename from which to read patterns for objects excluded or included from restore. The patterns are interpreted according to the same rules as `-n/--schema` for including objects in schemas, `-N/--exclude-schema` for excluding objects in schemas, `-P/--function` for restoring named functions, `-I/--index` for restoring named indexes, `-t/--table` for restoring named tables or `-T/--trigger` for restoring triggers. To read from `STDIN`, use `-` as the filename. The `--filter` option can be specified in conjunction with the above listed options for including or excluding objects, and can also be specified more than once for multiple filter files.

The file lists one database pattern per row, with the following format:

```
{ include | exclude } { function | index | schema | table | trigger } PATTERN
```

The first keyword specifies whether the objects matched by the pattern are to be included or excluded. The second keyword specifies the type of object to be filtered using the pattern:

- **function:** functions, works like the `-P/--function` option. This keyword can only be used with the **include** keyword.
- **index:** indexes, works like the `-I/--indexes` option. This keyword can only be used with the **include** keyword.
- **schema:** schemas, works like the `-n/--schema` and `-N/--exclude-schema` options.
- **table:** tables, works like the `-t/--table` option. This keyword can only be used with the **include** keyword.
- **trigger:** triggers, works like the `-T/--trigger` option. This keyword can only be used with the **include** keyword.

Lines starting with `#` are considered comments and ignored. Comments can be placed after an object pattern row as well. Blank lines are also ignored. See [Patterns](#) for how to perform quoting in patterns.

`-F format`
`--format=format`

Specify format of the archive. It is not necessary to specify the format, since `pg_restore` will determine the format automatically. If specified, it can be one of the following:

`c`
`custom`

The archive is in the custom format of `pg_dump`.

`d`
`directory`

The archive is a directory archive.

`t`
`tar`

The archive is a `tar` archive.

`-I index`
`--index=index`

Restore definition of named index only. Multiple indexes may be specified with multiple `-I` switches.

`-j number-of-jobs`
`--jobs=number-of-jobs`

Run the most time-consuming steps of `pg_restore` — those that load data, create indexes, or create constraints — concurrently, using up to *number-of-jobs* concurrent sessions. This option can dramatically reduce the time to restore a large database to a server running on a multiprocessor machine. This option is ignored when emitting a script rather than connecting directly to a database server.

Each job is one process or one thread, depending on the operating system, and uses a separate connection to the server.

The optimal value for this option depends on the hardware setup of the server, of the client, and of the network. Factors include the number of CPU cores and the disk setup. A good place to start is the number of CPU cores on the server, but values larger than that can also lead to faster restore times in many cases. Of course, values that are too high will lead to decreased performance because of thrashing.

Only the custom and directory archive formats are supported with this option. The input must be a regular file or directory (not, for example, a pipe or standard input). Also, multiple jobs cannot be used together with the option `--single-transaction`.

`-l`
`--list`

List the table of contents of the archive. The output of this operation can be used as input to the `-L` option. Note that if filtering switches such as `-n` or `-t` are used with `-l`, they will restrict the items listed.

`-L list-file`
`--use-list=list-file`

Restore only those archive elements that are listed in *list-file*, and restore them in the order they appear in the file. Note that if filtering switches such as `-n` or `-t` are used with `-L`, they will further restrict the items restored.

list-file is normally created by editing the output of a previous `-l` operation. Lines can be moved or removed, and can also be commented out by placing a semicolon (;) at the start of the line. See below for examples.

`-n schema`
`--schema=schema`

Restore only objects that are in the named schema. Multiple schemas may be specified with multiple `-n` switches. This can be combined with the `-t` option to restore just a specific table.

`-N schema`
`--exclude-schema=schema`

Do not restore objects that are in the named schema. Multiple schemas to be excluded may be specified with multiple `-N` switches.

When both `-n` and `-N` are given for the same schema name, the `-N` switch wins and the schema is excluded.

`-O`
`--no-owner`

Do not output commands to set ownership of objects to match the original database. By default, `pg_restore` issues `ALTER OWNER` or `SET SESSION AUTHORIZATION` statements to set ownership of created schema elements. These statements will fail unless the initial connection to the database is made by a superuser (or the same user that owns all of the objects in the script). With `-O`, any user name can be used for the initial connection, and this user will own all the created objects.

`-P function-name(argtype [, ...])`
`--function=function-name(argtype [, ...])`

Restore the named function only. Be careful to spell the function name and arguments exactly as they appear in the dump file's table of contents. Multiple functions may be specified with multiple `-P` switches.

`-R`
`--no-reconnect`

This option is obsolete but still accepted for backwards compatibility.

`-s`
`--schema-only`

Restore only the schema (data definitions), not data, to the extent that schema entries are present in the archive.

This option is the inverse of `--data-only`. It is similar to, but for historical reasons not identical to, specifying `--section=pre-data --section=post-data`.

(Do not confuse this with the `--schema` option, which uses the word “schema” in a different meaning.)

`-S username`
`--superuser=username`

Specify the superuser user name to use when disabling triggers. This is relevant only if `--disable-triggers` is used.

`-t table`
`--table=table`

Restore definition and/or data of only the named table. For this purpose, “table” includes views, materialized views, sequences, and foreign tables. Multiple tables can be selected by writing multiple `-t` switches. This option can be combined with the `-n` option to specify table(s) in a particular schema.

Note

When `-t` is specified, `pg_restore` makes no attempt to restore any other database objects that the selected table(s) might depend upon. Therefore, there is no guarantee that a specific-table restore into a clean database will succeed.

Note

This flag does not behave identically to the `-t` flag of `pg_dump`. There is not currently any provision for wild-card matching in `pg_restore`, nor can you include a schema name within its `-t`. And, while `pg_dump`'s `-t` flag will also dump subsidiary objects (such as indexes) of the selected table(s), `pg_restore`'s `-t` flag does not include such subsidiary objects.

Note

In versions prior to Postgres Pro 9.6, this flag matched only tables, not any other type of relation.

`-T trigger`
`--trigger=trigger`

Restore named trigger only. Multiple triggers may be specified with multiple `-T` switches.

`-v`
`--verbose`

Specifies verbose mode. This will cause `pg_restore` to output detailed object comments and start/stop times to the output file, and progress messages to standard error. Repeating the option causes additional debug-level messages to appear on standard error.

`-V`
`--version`

Print the `pg_restore` version and exit.

`-x`
`--no-privileges`
`--no-acl`

Prevent restoration of access privileges (grant/revoke commands).

`-1`
`--single-transaction`

Execute the restore as a single transaction (that is, wrap the emitted commands in `BEGIN/COMMIT`). This ensures that either all the commands complete successfully, or no changes are applied. This option implies `--exit-on-error`.

`--disable-triggers`

This option is relevant only when performing a data-only restore. It instructs `pg_restore` to execute commands to temporarily disable triggers on the target tables while the data is restored. Use this if you have referential integrity checks or other triggers on the tables that you do not want to invoke during data restore.

Presently, the commands emitted for `--disable-triggers` must be done as superuser. So you should also specify a superuser name with `-s` or, preferably, run `pg_restore` as a Postgres Pro superuser.

`--enable-row-security`

This option is relevant only when restoring the contents of a table which has row security. By default, `pg_restore` will set [row_security](#) to off, to ensure that all data is restored in to the table. If the user does not have sufficient privileges to bypass row security, then an error is thrown. This parameter instructs `pg_restore` to set [row_security](#) to on instead, allowing the user to attempt to restore the

contents of the table with row security enabled. This might still fail if the user does not have the right to insert the rows from the dump into the table.

Note that this option currently also requires the dump be in `INSERT` format, as `COPY FROM` does not support row security.

`--if-exists`

Use `DROP ... IF EXISTS` commands to drop objects in `--clean` mode. This suppresses “does not exist” errors that might otherwise be reported. This option is not valid unless `--clean` is also specified.

`--no-comments`

Do not output commands to restore comments, even if the archive contains them.

`--no-data-for-failed-tables`

By default, table data is restored even if the creation command for the table failed (e.g., because it already exists). With this option, data for such a table is skipped. This behavior is useful if the target database already contains the desired table contents. For example, auxiliary tables for Postgres Pro extensions such as PostGIS might already be loaded in the target database; specifying this option prevents duplicate or obsolete data from being loaded into them.

This option is effective only when restoring directly into a database, not when producing SQL script output.

`--no-publications`

Do not output commands to restore publications, even if the archive contains them.

`--no-security-labels`

Do not output commands to restore security labels, even if the archive contains them.

`--no-subscriptions`

Do not output commands to restore subscriptions, even if the archive contains them.

`--no-table-access-method`

Do not output commands to select table access methods. With this option, all objects will be created with whichever table access method is the default during restore.

`--no-tablespaces`

Do not output commands to select tablespaces. With this option, all objects will be created in whichever tablespace is the default during restore.

`--section=sectionname`

Only restore the named section. The section name can be `pre-data`, `data`, or `post-data`. This option can be specified more than once to select multiple sections. The default is to restore all sections.

The data section contains actual table data as well as large-object definitions. Post-data items consist of definitions of indexes, triggers, rules and constraints other than validated check constraints. Pre-data items consist of all other data definition items.

`--strict-names`

Require that each schema (`-n/--schema`) and table (`-t/--table`) qualifier match at least one schema/table in the backup file.

`--transaction-size=N`

Execute the restore as a series of transactions, each processing up to *N* database objects. This option implies `--exit-on-error`.

`--transaction-size` offers an intermediate choice between the default behavior (one transaction per SQL command) and `-1/--single-transaction` (one transaction for all restored objects). While `--single-transaction` has the least overhead, it may be impractical for large databases because the transaction will take a lock on each restored object, possibly exhausting the server's lock table space. Using `--transaction-size` with a size of a few thousand objects offers nearly the same performance benefits while capping the amount of lock table space needed.

`--use-set-session-authorization`

Output SQL-standard `SET SESSION AUTHORIZATION` commands instead of `ALTER OWNER` commands to determine object ownership. This makes the dump more standards-compatible, but depending on the history of the objects in the dump, might not restore properly.

`-?`

`--help`

Show help about `pg_restore` command line arguments, and exit.

`pg_restore` also accepts the following command line arguments for connection parameters:

`-h host`

`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket. The default is taken from the `PGHOST` environment variable, if set, else a Unix domain socket connection is attempted.

`-p port`

`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections. Defaults to the `PGPORT` environment variable, if set, or a compiled-in default.

`-U username`

`--username=username`

User name to connect as.

`-w`

`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`

`--password`

Force `pg_restore` to prompt for a password before connecting to a database.

This option is never essential, since `pg_restore` will automatically prompt for a password if the server demands password authentication. However, `pg_restore` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`--role=rolename`

Specifies a role name to be used to perform the restore. This option causes `pg_restore` to issue a `SET ROLE rolename` command after connecting to the database. It is useful when the authenticated user (specified by `-U`) lacks privileges needed by `pg_restore`, but can switch to a role with the required rights. Some installations have a policy against logging in directly as a superuser, and use of this option allows restores to be performed without violating the policy.

Environment

PGHOST
PGOPTIONS
PGPORT
PGUSER

Default connection parameters

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)). However, it does not read `PGDATABASE` when a database name is not supplied.

Diagnostics

When a direct database connection is specified using the `-d` option, `pg_restore` internally executes SQL statements. If you have problems running `pg_restore`, make sure you are able to select information from the database using, for example, `psql`. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Notes

If your installation has any local additions to the `template1` database, be careful to load the output of `pg_restore` into a truly empty database; otherwise you are likely to get errors due to duplicate definitions of the added objects. To make an empty database without any local additions, copy from `template0` not `template1`, for example:

```
CREATE DATABASE foo WITH TEMPLATE template0;
```

The limitations of `pg_restore` are detailed below.

- When restoring data to a pre-existing table and the option `--disable-triggers` is used, `pg_restore` emits commands to disable triggers on user tables before inserting the data, then emits commands to re-enable them after the data has been inserted. If the restore is stopped in the middle, the system catalogs might be left in the wrong state.
- `pg_restore` cannot restore large objects selectively; for instance, only those for a specific table. If an archive contains large objects, then all large objects will be restored, or none of them if they are excluded via `-L`, `-t`, or other options.

See also the [pg_dump](#) documentation for details on limitations of `pg_dump`.

Once restored, it is wise to run `ANALYZE` on each restored table so the optimizer has useful statistics; see [Section 24.1.3](#) and [Section 24.1.6](#) for more information.

Examples

Assume we have dumped a database called `mydb` into a custom-format dump file:

```
$ pg_dump -Fc mydb > db.dump
```

To drop the database and recreate it from the dump:

```
$ dropdb mydb
$ pg_restore -C -d postgres db.dump
```

The database named in the `-d` switch can be any database existing in the cluster; `pg_restore` only uses it to issue the `CREATE DATABASE` command for `mydb`. With `-C`, data is always restored into the database name that appears in the dump file.

To restore the dump into a new database called `newdb`:

```
$ createdb -T template0 newdb
$ pg_restore -d newdb db.dump
```

Notice we don't use `-C`, and instead connect directly to the database to be restored into. Also note that we clone the new database from `template0` not `template1`, to ensure it is initially empty.

To reorder database items, it is first necessary to dump the table of contents of the archive:

```
$ pg_restore -l db.dump > db.list
```

The listing file consists of a header and one line for each item, e.g.:

```
;
; Archive created at Mon Sep 14 13:55:39 2009
;   dbname: DBDEMOS
;   TOC Entries: 81
;   Compression: 9
;   Dump Version: 1.10-0
;   Format: CUSTOM
;   Integer: 4 bytes
;   Offset: 8 bytes
;   Dumped from database version: 8.3.5
;   Dumped by pg_dump version: 8.3.8
;
;
; Selected TOC Entries:
;
3; 2615 2200 SCHEMA - public pasha
1861; 0 0 COMMENT - SCHEMA public pasha
1862; 0 0 ACL - public pasha
317; 1247 17715 TYPE public composite pasha
319; 1247 25899 DOMAIN public domain0 pasha
```

Semicolons start a comment, and the numbers at the start of lines refer to the internal archive ID assigned to each item.

Lines in the file can be commented out, deleted, and reordered. For example:

```
10; 145433 TABLE map_resolutions postgres
;2; 145344 TABLE species postgres
;4; 145359 TABLE nt_header postgres
6; 145402 TABLE species_records postgres
;8; 145416 TABLE ss_old postgres
```

could be used as input to `pg_restore` and would only restore items 10 and 6, in that order:

```
$ pg_restore -L db.list db.dump
```

See Also

[pg_dump](#), [pg_dumpall](#), [psql](#)

pg-wrapper

pg-wrapper — manage Postgres Pro symbolic links

Synopsis

```
pg-wrapper links { update | remove }
```

Description

pg-wrapper is a shell script provided in the Postgres Pro distribution to manage Postgres Pro symbolic links for the provided programs and adjust the handling of SQL man pages on Linux systems. Since Postgres Pro is installed into `/opt/pgpro/sdm-17`, this is required to make client and server programs available in the standard system `PATH` and find the new SQL man pages. This setup is not performed automatically at installation time to avoid possible conflicts with other PostgreSQL-based products you may have installed, unless you have opted for using the `postgrespro-sdm-17` quick-install package.

pg-wrapper is provided as part of the `postgrespro-sdm-17-client` package. Once this package is installed, you can find pg-wrapper in the `install-dir/bin` directory, where `install-dir` is `/opt/pgpro/sdm-17`.

pg-wrapper must be run as root.

Options

pg-wrapper accepts the following command-line arguments:

```
links {update | remove}
```

Manage Postgres Pro symbolic links and SQL man pages:

- `update` — create or replace symbolic links for server and client programs provided with Postgres Pro, as well as add SQL man pages to the man page configuration file.
- `remove` — remove symbolic links and SQL man page support for the current Postgres Pro version.

Notes

If you are installing Postgres Pro from the `postgrespro-sdm-17` package, pg-wrapper is invoked automatically.

If you are installing individual Postgres Pro packages, you can run this script manually to create symbolic links to the provided client and server programs, as well as add SQL man pages to the man page configuration file.

For parallel installations with other PostgreSQL-based products, pg-wrapper behavior depends on whether the `update-alternatives` utility is supported by your system and the already installed PostgreSQL-based product:

- If `update-alternatives` is supported, pg-wrapper adds symbolic links to Postgres Pro programs into the alternative system, in accordance with their priority, as well as adds SQL man pages to the man page configuration file. For details on how to change the alternatives system priorities, see the man page for `update-alternatives` on your system.

If SQL man pages are already installed from a different product, the pages unique for the new installation will be displayed, while the previous installation will keep its priority for all the coinciding man pages. For the new SQL documentation to be displayed for all pages, you have to modify the system configuration, for example, change the `MANPATH` value.

- If `update-alternatives` is not supported, pg-wrapper updates the system configuration only if there are no conflicts with any programs or man pages already installed. Otherwise, pg-wrapper does not

create or update any program links and skips SQL man page integration. In this case, you can either continue using the already available program versions and SQL man pages, or modify the `PATH` and `MANPATH` settings manually.

For details on binary installation specifics on Linux, see [Section 17.1](#).

pg_verifybackup

pg_verifybackup — verify the integrity of a base backup of a Postgres Pro cluster

Synopsis

```
pg_verifybackup [option...]
```

Description

pg_verifybackup is used to check the integrity of a database cluster backup taken using pg_basebackup against a backup_manifest generated by the server at the time of the backup. The backup must be stored in the "plain" format; a "tar" format backup can be checked after extracting it.

It is important to note that the validation which is performed by pg_verifybackup does not and cannot include every check which will be performed by a running server when attempting to make use of the backup. Even if you use this tool, you should still perform test restores and verify that the resulting databases work as expected and that they appear to contain the correct data. However, pg_verifybackup can detect many problems that commonly occur due to storage problems or user error.

Backup verification proceeds in four stages. First, pg_verifybackup reads the backup_manifest file. If that file does not exist, cannot be read, is malformed, fails to match the system identifier with pg_control of the backup directory or fails verification against its own internal checksum, pg_verifybackup will terminate with a fatal error.

Second, pg_verifybackup will attempt to verify that the data files currently stored on disk are exactly the same as the data files which the server intended to send, with some exceptions that are described below. Extra and missing files will be detected, with a few exceptions. This step will ignore the presence or absence of, or any modifications to, postgresql.auto.conf, standby.signal, and recovery.signal, because it is expected that these files may have been created or modified as part of the process of taking the backup. It also won't complain about a backup_manifest file in the target directory or about anything inside pg_wal, even though these files won't be listed in the backup manifest. Only files are checked; the presence or absence of directories is not verified, except indirectly: if a directory is missing, any files it should have contained will necessarily also be missing.

Next, pg_verifybackup will checksum all the files, compare the checksums against the values in the manifest, and emit errors for any files for which the computed checksum does not match the checksum stored in the manifest. This step is not performed for any files which produced errors in the previous step, since they are already known to have problems. Files which were ignored in the previous step are also ignored in this step.

Finally, pg_verifybackup will use the manifest to verify that the write-ahead log records which will be needed to recover the backup are present and that they can be read and parsed. The backup_manifest contains information about which write-ahead log records will be needed, and pg_verifybackup will use that information to invoke pg_waldump to parse those write-ahead log records. The --quiet flag will be used, so that pg_waldump will only report errors, without producing any other output. While this level of verification is sufficient to detect obvious problems such as a missing file or one whose internal checksums do not match, they aren't extensive enough to detect every possible problem that might occur when attempting to recover. For instance, a server bug that produces write-ahead log records that have the correct checksums but specify nonsensical actions can't be detected by this method.

Note that if extra WAL files which are not required to recover the backup are present, they will not be checked by this tool, although a separate invocation of pg_waldump could be used for that purpose. Also note that WAL verification is version-specific: you must use the version of pg_verifybackup, and thus of pg_waldump, which pertains to the backup being checked. In contrast, the data file integrity checks should work with any version of the server that generates a backup_manifest file.

Options

pg_verifybackup accepts the following command-line arguments:

`-e`
`--exit-on-error`

Exit as soon as a problem with the backup is detected. If this option is not specified, `pg_verifybackup` will continue checking the backup even after a problem has been detected, and will report all problems detected as errors.

`-i path`
`--ignore=path`

Ignore the specified file or directory, which should be expressed as a relative path name, when comparing the list of data files actually present in the backup to those listed in the `backup_manifest` file. If a directory is specified, this option affects the entire subtree rooted at that location. Complaints about extra files, missing files, file size differences, or checksum mismatches will be suppressed if the relative path name matches the specified path name. This option can be specified multiple times.

`-m path`
`--manifest-path=path`

Use the manifest file at the specified path, rather than one located in the root of the backup directory.

`-n`
`--no-parse-wal`

Don't attempt to parse write-ahead log data that will be needed to recover from this backup.

`-P`
`--progress`

Enable progress reporting. Turning this on will deliver a progress report while verifying checksums.

This option cannot be used together with the option `--quiet`.

`-q`
`--quiet`

Don't print anything when a backup is successfully verified.

`-s`
`--skip-checksums`

Do not verify data file checksums. The presence or absence of files and the sizes of those files will still be checked. This is much faster, because the files themselves do not need to be read.

`-w path`
`--wal-directory=path`

Try to parse WAL files stored in the specified directory, rather than in `pg_wal`. This may be useful if the backup is stored in a separate location from the WAL archive.

Other options are also available:

`-V`
`--version`

Print the `pg_verifybackup` version and exit.

`-?`
`--help`

Show help about `pg_verifybackup` command line arguments, and exit.

Examples

To create a base backup of the server at `mydbserver` and verify the integrity of the backup:

```
$ pg_basebackup -h mydbserver -D /usr/local/pgsql/data
$ pg_verifybackup /usr/local/pgsql/data
```

To create a base backup of the server at `mydbserver`, move the manifest somewhere outside the backup directory, and verify the backup:

```
$ pg_basebackup -h mydbserver -D /usr/local/pgsql/backup1234
$ mv /usr/local/pgsql/backup1234/backup_manifest /my/secure/location/
  backup_manifest.1234
$ pg_verifybackup -m /my/secure/location/backup_manifest.1234 /usr/local/pgsql/
  backup1234
```

To verify a backup while ignoring a file that was added manually to the backup directory, and also skipping checksum verification:

```
$ pg_basebackup -h mydbserver -D /usr/local/pgsql/data
$ edit /usr/local/pgsql/data/note.to.self
$ pg_verifybackup --ignore=note.to.self --skip-checksums /usr/local/pgsql/data
```

psql

psql — Postgres Pro interactive terminal

Synopsis

```
psql [option...] [dbname [username]]
```

Description

psql is a terminal-based front-end to Postgres Pro. It enables you to type in queries interactively, issue them to Postgres Pro, and see the query results. Alternatively, input can be from a file or from command line arguments. In addition, psql provides a number of meta-commands and various shell-like features to facilitate writing scripts and automating a wide variety of tasks.

Options

```
-a  
--echo-all
```

Print all nonempty input lines to standard output as they are read. (This does not apply to lines read interactively.) This is equivalent to setting the variable `ECHO` to `all`.

```
-A  
--no-align
```

Switches to unaligned output mode. (The default output mode is `aligned`.) This is equivalent to `\pset format unaligned`.

```
-b  
--echo-errors
```

Print failed SQL commands to standard error output. This is equivalent to setting the variable `ECHO` to `errors`.

```
-c command  
--command=command
```

Specifies that psql is to execute the given command string, *command*. This option can be repeated and combined in any order with the `-f` option. When either `-c` or `-f` is specified, psql does not read commands from standard input; instead it terminates after processing all the `-c` and `-f` options in sequence.

command must be either a command string that is completely parsable by the server (i.e., it contains no psql-specific features), or a single backslash command. Thus you cannot mix SQL and psql meta-commands within a `-c` option. To achieve that, you could use repeated `-c` options or pipe the string into psql, for example:

```
psql -c '\x' -c 'SELECT * FROM foo;'
```

or

```
echo '\x \ SELECT * FROM foo;' | psql
```

(`\` is the separator meta-command.)

Each SQL command string passed to `-c` is sent to the server as a single request. Because of this, the server executes it as a single transaction even if the string contains multiple SQL commands, unless there are explicit `BEGIN/COMMIT` commands included in the string to divide it into multiple transactions. (See [Section 56.2.2.1](#) for more details about how the server handles multi-query strings.)

If having several commands executed in one transaction is not desired, use repeated `-c` commands or feed multiple commands to psql's standard input, either using echo as illustrated above, or via a shell here-document, for example:

```
psql <<EOF
\x
SELECT * FROM foo;
EOF
```

`--csv`

Switches to CSV (Comma-Separated Values) output mode. This is equivalent to `\pset format csv`.

`-d dbname`

`--dbname=dbname`

Specifies the name of the database to connect to. This is equivalent to specifying `dbname` as the first non-option argument on the command line. The `dbname` can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

`-e`

`--echo-queries`

Copy all SQL commands sent to the server to standard output as well. This is equivalent to setting the variable `ECHO` to `queries`.

`-E`

`--echo-hidden`

Echo the actual queries generated by `\d` and other backslash commands. You can use this to study psql's internal operations. This is equivalent to setting the variable `ECHO_HIDDEN` to `on`.

`-f filename`

`--file=filename`

Read commands from the file `filename`, rather than standard input. This option can be repeated and combined in any order with the `-c` option. When either `-c` or `-f` is specified, psql does not read commands from standard input; instead it terminates after processing all the `-c` and `-f` options in sequence. Except for that, this option is largely equivalent to the meta-command `\i`.

If `filename` is `-` (hyphen), then standard input is read until an EOF indication or `\q` meta-command. This can be used to intersperse interactive input with input from files. Note however that Readline is not used in this case (much as if `-n` had been specified).

Using this option is subtly different from writing `psql < filename`. In general, both will do what you expect, but using `-f` enables some nice features such as error messages with line numbers. There is also a slight chance that using this option will reduce the start-up overhead. On the other hand, the variant using the shell's input redirection is (in theory) guaranteed to yield exactly the same output you would have received had you entered everything by hand.

`-F separator`

`--field-separator=separator`

Use `separator` as the field separator for unaligned output. This is equivalent to `\pset fieldsep` or `\f`.

`-h hostname`

`--host=hostname`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix-domain socket.

-H
--html

Switches to HTML output mode. This is equivalent to `\pset format html` or the `\H` command.

-l
--list

List all available databases, then exit. Other non-connection options are ignored. This is similar to the meta-command `\list`.

When this option is used, `psql` will connect to the database `postgres`, unless a different database is named on the command line (option `-d` or non-option argument, possibly via a service entry, but not via an environment variable).

-L *filename*
--log-file=*filename*

Write all query output into file *filename*, in addition to the normal output destination.

-n
--no-readline

Do not use Readline for line editing and do not use the command history (see [the section called “Command-Line Editing”](#) below).

-o *filename*
--output=*filename*

Put all query output into file *filename*. This is equivalent to the command `\o`.

-p *port*
--port=*port*

Specifies the TCP port or the local Unix-domain socket file extension on which the server is listening for connections. Defaults to the value of the `PGPORT` environment variable or, if not set, to the port specified at compile time, usually 5432.

-P *assignment*
--pset=*assignment*

Specifies printing options, in the style of `\pset`. Note that here you have to separate name and value with an equal sign instead of a space. For example, to set the output format to LaTeX, you could write `-P format=latex`.

-q
--quiet

Specifies that `psql` should do its work quietly. By default, it prints welcome messages and various informational output. If this option is used, none of this happens. This is useful with the `-c` option. This is equivalent to setting the variable `QUIET` to `on`.

-R *separator*
--record-separator=*separator*

Use *separator* as the record separator for unaligned output. This is equivalent to `\pset recordsep`.

-s
--single-step

Run in single-step mode. That means the user is prompted before each command is sent to the server, with the option to cancel execution as well. Use this to debug scripts.

-S
--single-line

Runs in single-line mode where a newline terminates an SQL command, as a semicolon does.

Note

This mode is provided for those who insist on it, but you are not necessarily encouraged to use it. In particular, if you mix SQL and meta-commands on a line the order of execution might not always be clear to the inexperienced user.

-t
--tuples-only

Turn off printing of column names and result row count footers, etc. This is equivalent to `\t` or `\pset tuples_only`.

-T *table_options*
--table-attr=*table_options*

Specifies options to be placed within the HTML `table` tag. See `\pset tableattr` for details.

-U *username*
--username=*username*

Connect to the database as the user *username* instead of the default. (You must have permission to do so, of course.)

-v *assignment*
--set=*assignment*
--variable=*assignment*

Perform a variable assignment, like the `\set` meta-command. Note that you must separate name and value, if any, by an equal sign on the command line. To unset a variable, leave off the equal sign. To set a variable with an empty value, use the equal sign but leave off the value. These assignments are done during command line processing, so variables that reflect connection state will get overwritten later.

-V
--version

Print the psql version and exit.

-w
--no-password

Never issue a password prompt. If the server requires password authentication and a password is not available from other sources such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

Note that this option will remain set for the entire session, and so it affects uses of the meta-command `\connect` as well as the initial connection attempt.

-W
--password

Force psql to prompt for a password before connecting to a database, even if the password will not be used.

If the server requires password authentication and a password is not available from other sources such as a `.pgpass` file, psql will prompt for a password in any case. However, psql will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-W` to avoid the extra connection attempt.

Note that this option will remain set for the entire session, and so it affects uses of the meta-command `\connect` as well as the initial connection attempt.

`-x`
`--expanded`

Turn on the expanded table formatting mode. This is equivalent to `\x` or `\pset expanded`.

`-X`
`--no-psqlrc`

Do not read the start-up file (neither the system-wide `psqlrc` file nor the user's `~/.psqlrc` file).

`-z`
`--field-separator-zero`

Set the field separator for unaligned output to a zero byte. This is equivalent to `\pset fieldsep_zero`.

`-0`
`--record-separator-zero`

Set the record separator for unaligned output to a zero byte. This is useful for interfacing, for example, with `xargs -0`. This is equivalent to `\pset recordsep_zero`.

`-1`
`--single-transaction`

This option can only be used in combination with one or more `-c` and/or `-f` options. It causes `psql` to issue a `BEGIN` command before the first such option and a `COMMIT` command after the last one, thereby wrapping all the commands into a single transaction. If any of the commands fails and the variable `ON_ERROR_STOP` was set, a `ROLLBACK` command is sent instead. This ensures that either all the commands complete successfully, or no changes are applied.

If the commands themselves contain `BEGIN`, `COMMIT`, or `ROLLBACK`, this option will not have the desired effects. Also, if an individual command cannot be executed inside a transaction block, specifying this option will cause the whole transaction to fail.

`-?`
`--help[=topic]`

Show help about `psql` and exit. The optional *topic* parameter (defaulting to `options`) selects which part of `psql` is explained: `commands` describes `psql`'s backslash commands; `options` describes the command-line options that can be passed to `psql`; and `variables` shows help about `psql` configuration variables.

Exit Status

`psql` returns 0 to the shell if it finished normally, 1 if a fatal error of its own occurs (e.g., out of memory, file not found), 2 if the connection to the server went bad and the session was not interactive, and 3 if an error occurred in a script and the variable `ON_ERROR_STOP` was set.

Usage

Connecting to a Database

`psql` is a regular Postgres Pro client application. In order to connect to a database you need to know the name of your target database, the host name and port number of the server, and what database user name you want to connect as. `psql` can be told about those parameters via command line options, namely `-d`, `-h`, `-p`, and `-U` respectively. If an argument is found that does not belong to any option it will be interpreted as the database name (or the database user name, if the database name is already given). Not all of these options are required; there are useful defaults. If you omit the host name, `psql` will connect via a Unix-domain socket to a server on the local host, or via TCP/IP to `localhost` on Windows. The default

port number is determined at compile time. Since the database server uses the same default, you will not have to specify the port in most cases. The default database user name is your operating-system user name. Once the database user name is determined, it is used as the default database name. Note that you cannot just connect to any database under any database user name. Your database administrator should have informed you about your access rights.

When the defaults aren't quite right, you can save yourself some typing by setting the environment variables `PGDATABASE`, `PGHOST`, `PGPORT` and/or `PGUSER` to appropriate values. (For additional environment variables, see [Section 35.15](#).) It is also convenient to have a `~/.pgpass` file to avoid regularly having to type in passwords. See [Section 35.16](#) for more information.

An alternative way to specify connection parameters is in a *conninfo* string or a URI, which is used instead of a database name. This mechanism give you very wide control over the connection. For example:

```
$ psql "service=myservice sslmode=require"
$ psql postgresql://dbmaster:5433/mydb?sslmode=require
```

This way you can also use LDAP for connection parameter lookup as described in [Section 35.18](#). See [Section 35.1.2](#) for more information on all the available connection options.

If the connection could not be made for any reason (e.g., insufficient privileges, server is not running on the targeted host, etc.), `psql` will return an error and terminate.

If both standard input and standard output are a terminal, then `psql` sets the client encoding to “auto”, which will detect the appropriate client encoding from the locale settings (`LC_CTYPE` environment variable on Unix systems). If this doesn't work out as expected, the client encoding can be overridden using the environment variable `PGCLIENTENCODING`.

Entering SQL Commands

In normal operation, `psql` provides a prompt with the name of the database to which `psql` is currently connected, followed by the string `=>`. For example:

```
$ psql testdb
psql (17.5.1)
Type "help" for help.

testdb=>
```

At the prompt, the user can type in SQL commands. Ordinarily, input lines are sent to the server when a command-terminating semicolon is reached. An end of line does not terminate a command. Thus commands can be spread over several lines for clarity. If the command was sent and executed without error, the results of the command are displayed on the screen.

If untrusted users have access to a database that has not adopted a [secure schema usage pattern](#), begin your session by removing publicly-writable schemas from `search_path`. One can add `options=-csearch_path=` to the connection string or issue `SELECT pg_catalog.set_config('search_path', '', false)` before other SQL commands. This consideration is not specific to `psql`; it applies to every interface for executing arbitrary SQL commands.

Whenever a command is executed, `psql` also polls for asynchronous notification events generated by [LISTEN](#) and [NOTIFY](#).

While C-style block comments are passed to the server for processing and removal, SQL-standard comments are removed by `psql`.

Meta-Commands

Anything you enter in `psql` that begins with an unquoted backslash is a `psql` meta-command that is processed by `psql` itself. These commands make `psql` more useful for administration or scripting. Meta-commands are often called slash or backslash commands.

The format of a psql command is the backslash, followed immediately by a command verb, then any arguments. The arguments are separated from the command verb and each other by any number of whitespace characters.

To include whitespace in an argument you can quote it with single quotes. To include a single quote in an argument, write two single quotes within single-quoted text. Anything contained in single quotes is furthermore subject to C-like substitutions for `\n` (new line), `\t` (tab), `\b` (backspace), `\r` (carriage return), `\f` (form feed), `\digits` (octal), and `\xdigits` (hexadecimal). A backslash preceding any other character within single-quoted text quotes that single character, whatever it is.

If an unquoted colon (`:`) followed by a psql variable name appears within an argument, it is replaced by the variable's value, as described in [SQL Interpolation](#) below. The forms `:'variable_name'` and `:"variable_name"` described there work as well. The `: {?variable_name}` syntax allows testing whether a variable is defined. It is substituted by `TRUE` or `FALSE`. Escaping the colon with a backslash protects it from substitution.

Within an argument, text that is enclosed in backquotes (```) is taken as a command line that is passed to the shell. The output of the command (with any trailing newline removed) replaces the backquoted text. Within the text enclosed in backquotes, no special quoting or other processing occurs, except that appearances of `:variable_name` where `variable_name` is a psql variable name are replaced by the variable's value. Also, appearances of `:'variable_name'` are replaced by the variable's value suitably quoted to become a single shell command argument. (The latter form is almost always preferable, unless you are very sure of what is in the variable.) Because carriage return and line feed characters cannot be safely quoted on all platforms, the `:'variable_name'` form prints an error message and does not substitute the variable value when such characters appear in the value.

Some commands take an SQL identifier (such as a table name) as argument. These arguments follow the syntax rules of SQL: Unquoted letters are forced to lowercase, while double quotes (`"`) protect letters from case conversion and allow incorporation of whitespace into the identifier. Within double quotes, paired double quotes reduce to a single double quote in the resulting name. For example, `FOO"BAR"BAZ` is interpreted as `fooBARbaz`, and `"A weird" name"` becomes `A weird name`.

Parsing for arguments stops at the end of the line, or when another unquoted backslash is found. An unquoted backslash is taken as the beginning of a new meta-command. The special sequence `\\` (two backslashes) marks the end of arguments and continues parsing SQL commands, if any. That way SQL and psql commands can be freely mixed on a line. But in any case, the arguments of a meta-command cannot continue beyond the end of the line.

Many of the meta-commands act on the *current query buffer*. This is simply a buffer holding whatever SQL command text has been typed but not yet sent to the server for execution. This will include previous input lines as well as any text appearing before the meta-command on the same line.

The following meta-commands are defined:

`\a`

If the current table output format is unaligned, it is switched to aligned. If it is not unaligned, it is set to unaligned. This command is kept for backwards compatibility. See `\pset` for a more general solution.

`\bind [parameter] ...`

Sets query parameters for the next query execution, with the specified parameters passed for any parameter placeholders (`$1` etc.).

Example:

```
INSERT INTO tbl1 VALUES ($1, $2) \bind 'first value' 'second value' \g
```

This also works for query-execution commands besides `\g`, such as `\gx` and `\gset`.

This command causes the extended query protocol (see [Section 56.1.2](#)) to be used, unlike normal psql operation, which uses the simple query protocol. So this command can be useful to test the

extended query protocol from psql. (The extended query protocol is used even if the query has no parameters and this command specifies zero parameters.) This command affects only the next query executed; all subsequent queries will use the simple query protocol by default.

```
\c or \connect [ -reuse-previous=on/off ] [ dbname [ username ] [ host ] [ port ] |
conninfo ]
```

Establishes a new connection to a Postgres Pro server. The connection parameters to use can be specified either using a positional syntax (one or more of database name, user, host, and port), or using a *conninfo* connection string as detailed in [Section 35.1.1](#). If no arguments are given, a new connection is made using the same parameters as before.

Specifying any of *dbname*, *username*, *host* or *port* as `-` is equivalent to omitting that parameter.

The new connection can re-use connection parameters from the previous connection; not only database name, user, host, and port, but other settings such as *sslmode*. By default, parameters are re-used in the positional syntax, but not when a *conninfo* string is given. Passing a first argument of `-reuse-previous=on` or `-reuse-previous=off` overrides that default. If parameters are re-used, then any parameter not explicitly specified as a positional parameter or in the *conninfo* string is taken from the existing connection's parameters. An exception is that if the *host* setting is changed from its previous value using the positional syntax, any *hostaddr* setting present in the existing connection's parameters is dropped. Also, any password used for the existing connection will be re-used only if the user, host, and port settings are not changed. When the command neither specifies nor reuses a particular parameter, the libpq default is used.

If the new connection is successfully made, the previous connection is closed. If the connection attempt fails (wrong user name, access denied, etc.), the previous connection will be kept if psql is in interactive mode. But when executing a non-interactive script, the old connection is closed and an error is reported. That may or may not terminate the script; if it does not, all database-accessing commands will fail until another `\connect` command is successfully executed. This distinction was chosen as a user convenience against typos on the one hand, and a safety mechanism that scripts are not accidentally acting on the wrong database on the other hand. Note that whenever a `\connect` command attempts to re-use parameters, the values re-used are those of the last successful connection, not of any failed attempts made subsequently. However, in the case of a non-interactive `\connect` failure, no parameters are allowed to be re-used later, since the script would likely be expecting the values from the failed `\connect` to be re-used.

Examples:

```
=> \c mydb myuser host.dom 6432
=> \c service=foo
=> \c "host=localhost port=5432 dbname=mydb connect_timeout=10 sslmode=disable"
=> \c -reuse-previous=on sslmode=require      -- changes only sslmode
=> \c postgresql://tom@localhost/mydb?application_name=myapp
```

```
\C [ title ]
```

Sets the title of any tables being printed as the result of a query or unset any such title. This command is equivalent to `\pset title title`. (The name of this command derives from “caption”, as it was previously only used to set the caption in an HTML table.)

```
\cd [ directory ]
```

Changes the current working directory to *directory*. Without argument, changes to the current user's home directory.

Tip

To print your current working directory, use `\! pwd`.

`\conninfo`

Outputs information about the current database connection.

```
\copy { table [ ( column_list ) ] } from { 'filename' | program 'command' | stdin |
pstdin } [ [ with ] ( option [, ...] ) ] [ where condition ]
\copy { table [ ( column_list ) ] | ( query ) } to { 'filename' | program 'command' |
stdout | pstdout } [ [ with ] ( option [, ...] ) ]
```

Performs a frontend (client) copy. This is an operation that runs an SQL `COPY` command, but instead of the server reading or writing the specified file, psql reads or writes the file and routes the data between the server and the local file system. This means that file accessibility and privileges are those of the local user, not the server, and no SQL superuser privileges are required.

When `program` is specified, `command` is executed by psql and the data passed from or to `command` is routed between the server and the client. Again, the execution privileges are those of the local user, not the server, and no SQL superuser privileges are required.

For `\copy ... from stdin`, data rows are read from the same source that issued the command, continuing until `\.` is read or the stream reaches EOF. This option is useful for populating tables in-line within an SQL script file. For `\copy ... to stdout`, output is sent to the same place as psql command output, and the `COPY count` command status is not printed (since it might be confused with a data row). To read/write psql's standard input or output regardless of the current command source or `\o` option, write `from pstdin` or `to pstdout`.

The syntax of this command is similar to that of the SQL `COPY` command. All options other than the data source/destination are as specified for `COPY`. Because of this, special parsing rules apply to the `\copy` meta-command. Unlike most other meta-commands, the entire remainder of the line is always taken to be the arguments of `\copy`, and neither variable interpolation nor backquote expansion are performed in the arguments.

Tip

Another way to obtain the same result as `\copy ... to` is to use the SQL `COPY ... TO STDOUT` command and terminate it with `\g filename` or `\g |program`. Unlike `\copy`, this method allows the command to span multiple lines; also, variable interpolation and backquote expansion can be used.

Tip

These operations are not as efficient as the SQL `COPY` command with a file or program data source or destination, because all data must pass through the client/server connection. For large amounts of data the SQL command might be preferable. Also, because of this pass-through method, `\copy ... from` in CSV mode will erroneously treat a `\.` data value alone on a line as an end-of-input marker.

`\copyright`

Shows the copyright and distribution terms of Postgres Pro.

```
\crosstabview [ colV [ colH [ colD [ sortcolH ] ] ] ]
```

Executes the current query buffer (like `\g`) and shows the results in a crosstab grid. The query must return at least three columns. The output column identified by `colV` becomes a vertical header and the output column identified by `colH` becomes a horizontal header. `colD` identifies the output column to display within the grid. `sortcolH` identifies an optional sort column for the horizontal header.

Each column specification can be a column number (starting at 1) or a column name. The usual SQL case folding and quoting rules apply to column names. If omitted, *colV* is taken as column 1 and *colH* as column 2. *colH* must differ from *colV*. If *colD* is not specified, then there must be exactly three columns in the query result, and the column that is neither *colV* nor *colH* is taken to be *colD*.

The vertical header, displayed as the leftmost column, contains the values found in column *colV*, in the same order as in the query results, but with duplicates removed.

The horizontal header, displayed as the first row, contains the values found in column *colH*, with duplicates removed. By default, these appear in the same order as in the query results. But if the optional *sortcolH* argument is given, it identifies a column whose values must be integer numbers, and the values from *colH* will appear in the horizontal header sorted according to the corresponding *sortcolH* values.

Inside the crosstab grid, for each distinct value *x* of *colH* and each distinct value *y* of *colV*, the cell located at the intersection (*x*, *y*) contains the value of the *colD* column in the query result row for which the value of *colH* is *x* and the value of *colV* is *y*. If there is no such row, the cell is empty. If there are multiple such rows, an error is reported.

`\d[S+] [pattern]`

For each relation (table, view, materialized view, index, sequence, or foreign table) or composite type matching the *pattern*, show all columns, their types, the tablespace (if not the default) and any special attributes such as `NOT NULL` or defaults. Associated indexes, constraints, rules, and triggers are also shown. For foreign tables, the associated foreign server is shown as well. (“Matching the pattern” is defined in [Patterns](#) below.)

For some types of relation, `\d` shows additional information for each column: column values for sequences, indexed expressions for indexes, and foreign data wrapper options for foreign tables.

The command form `\d+` is identical, except that more information is displayed: any comments associated with the columns of the table are shown, as is the presence of OIDs in the table, the view definition if the relation is a view, a non-default [replica identity](#) setting and the [access method](#) name if the relation has an access method.

By default, only user-created objects are shown; supply a *pattern* or the *s* modifier to include system objects.

Note

If `\d` is used without a *pattern* argument, it is equivalent to `\dtvmsE` which will show a list of all visible tables, views, materialized views, sequences and foreign tables. This is purely a convenience measure.

`\da[S] [pattern]`

Lists aggregate functions, together with their return type and the data types they operate on. If *pattern* is specified, only aggregates whose names match the pattern are shown. By default, only user-created objects are shown; supply a *pattern* or the *s* modifier to include system objects.

`\dA[+] [pattern]`

Lists access methods. If *pattern* is specified, only access methods whose names match the pattern are shown. If *+* is appended to the command name, each access method is listed with its associated handler function and description.

`\dAc[+] [access-method-pattern [input-type-pattern]`

Lists operator classes (see [Section 39.16.1](#)). If *access-method-pattern* is specified, only operator classes associated with access methods whose names match that pattern are listed. If *input-type-*

pattern is specified, only operator classes associated with input types whose names match that pattern are listed. If + is appended to the command name, each operator class is listed with its associated operator family and owner.

`\dAf[+] [access-method-pattern [input-type-pattern]]`

Lists operator families (see [Section 39.16.5](#)). If *access-method-pattern* is specified, only operator families associated with access methods whose names match that pattern are listed. If *input-type-pattern* is specified, only operator families associated with input types whose names match that pattern are listed. If + is appended to the command name, each operator family is listed with its owner.

`\dAo[+] [access-method-pattern [operator-family-pattern]]`

Lists operators associated with operator families (see [Section 39.16.2](#)). If *access-method-pattern* is specified, only members of operator families associated with access methods whose names match that pattern are listed. If *operator-family-pattern* is specified, only members of operator families whose names match that pattern are listed. If + is appended to the command name, each operator is listed with its sort operator family (if it is an ordering operator).

`\dAp[+] [access-method-pattern [operator-family-pattern]]`

Lists support functions associated with operator families (see [Section 39.16.3](#)). If *access-method-pattern* is specified, only functions of operator families associated with access methods whose names match that pattern are listed. If *operator-family-pattern* is specified, only functions of operator families whose names match that pattern are listed. If + is appended to the command name, functions are displayed verbosely, with their actual parameter lists.

`\db[+] [pattern]`

Lists tablespaces. If *pattern* is specified, only tablespaces whose names match the pattern are shown. If + is appended to the command name, each tablespace is listed with its associated options, on-disk size, permissions and description.

`\dc[S+] [pattern]`

Lists conversions between character-set encodings. If *pattern* is specified, only conversions whose names match the pattern are listed. By default, only user-created objects are shown; supply a pattern or the *s* modifier to include system objects. If + is appended to the command name, each object is listed with its associated description.

`\dconfig[+] [pattern]`

Lists server configuration parameters and their values. If *pattern* is specified, only parameters whose names match the pattern are listed. Without a *pattern*, only parameters that are set to non-default values are listed. (Use `\dconfig *` to see all parameters.) If + is appended to the command name, each parameter is listed with its data type, context in which the parameter can be set, and access privileges (if non-default access privileges have been granted).

`\dC[+] [pattern]`

Lists type casts. If *pattern* is specified, only casts whose source or target types match the pattern are listed. If + is appended to the command name, each object is listed with its associated description.

`\dd[S] [pattern]`

Shows the descriptions of objects of type `constraint`, `operator class`, `operator family`, `rule`, and `trigger`. All other comments may be viewed by the respective backslash commands for those object types.

`\dd` displays descriptions for objects matching the *pattern*, or of visible objects of the appropriate type if no argument is given. But in either case, only objects that have a description are listed. By default, only user-created objects are shown; supply a pattern or the *s* modifier to include system objects.

Descriptions for objects can be created with the [COMMENT](#) SQL command.

`\d[S+] [pattern]`

Lists domains. If *pattern* is specified, only domains whose names match the pattern are shown. By default, only user-created objects are shown; supply a pattern or the *s* modifier to include system objects. If *+* is appended to the command name, each object is listed with its associated permissions and description.

`\ddp [pattern]`

Lists default access privilege settings. An entry is shown for each role (and schema, if applicable) for which the default privilege settings have been changed from the built-in defaults. If *pattern* is specified, only entries whose role name or schema name matches the pattern are listed.

The [ALTER DEFAULT PRIVILEGES](#) command is used to set default access privileges. The meaning of the privilege display is explained in [Section 5.8](#).

`\dE[S+] [pattern]`

`\di[S+] [pattern]`

`\dm[S+] [pattern]`

`\ds[S+] [pattern]`

`\dt[S+] [pattern]`

`\dv[S+] [pattern]`

In this group of commands, the letters *E*, *i*, *m*, *s*, *t*, and *v* stand for foreign table, index, materialized view, sequence, table, and view, respectively. You can specify any or all of these letters, in any order, to obtain a listing of objects of these types. For example, `\dti` lists tables and indexes. If *+* is appended to the command name, each object is listed with its persistence status (permanent, temporary, or unlogged), physical size on disk, and associated description if any. If *pattern* is specified, only objects whose names match the pattern are listed. By default, only user-created objects are shown; supply a pattern or the *s* modifier to include system objects.

`\des[+] [pattern]`

Lists foreign servers (mnemonic: “external servers”). If *pattern* is specified, only those servers whose name matches the pattern are listed. If the form `\des+` is used, a full description of each server is shown, including the server's access privileges, type, version, options, and description.

`\det[+] [pattern]`

Lists foreign tables (mnemonic: “external tables”). If *pattern* is specified, only entries whose table name or schema name matches the pattern are listed. If the form `\det+` is used, generic options and the foreign table description are also displayed.

`\deu[+] [pattern]`

Lists user mappings (mnemonic: “external users”). If *pattern* is specified, only those mappings whose user names match the pattern are listed. If the form `\deu+` is used, additional information about each mapping is shown.

Caution

`\deu+` might also display the user name and password of the remote user, so care should be taken not to disclose them.

`\dew[+] [pattern]`

Lists foreign-data wrappers (mnemonic: “external wrappers”). If *pattern* is specified, only those foreign-data wrappers whose name matches the pattern are listed. If the form `\dew+` is used, the access privileges, options, and description of the foreign-data wrapper are also shown.

`\df[anptwS+] [pattern [arg_pattern ...]]`

Lists functions, together with their result data types, argument data types, and function types, which are classified as “agg” (aggregate), “normal”, “procedure”, “trigger”, or “window”. To display only functions of specific type(s), add the corresponding letters *a*, *n*, *p*, *t*, or *w* to the command. If *pattern* is specified, only functions whose names match the pattern are shown. Any additional arguments are type-name patterns, which are matched to the type names of the first, second, and so on arguments of the function. (Matching functions can have more arguments than what you specify. To prevent that, write a dash - as the last *arg_pattern*.) By default, only user-created objects are shown; supply a pattern or the *S* modifier to include system objects. If the form `\df+` is used, additional information about each function is shown, including volatility, parallel safety, owner, security classification, access privileges, language, internal name (for C and internal functions only), and description. Source code for a specific function can be seen using `\sf`.

`\dF[+] [pattern]`

Lists text search configurations. If *pattern* is specified, only configurations whose names match the pattern are shown. If the form `\dF+` is used, a full description of each configuration is shown, including the underlying text search parser and the dictionary list for each parser token type.

`\dFd[+] [pattern]`

Lists text search dictionaries. If *pattern* is specified, only dictionaries whose names match the pattern are shown. If the form `\dFd+` is used, additional information is shown about each selected dictionary, including the underlying text search template and the option values.

`\dFp[+] [pattern]`

Lists text search parsers. If *pattern* is specified, only parsers whose names match the pattern are shown. If the form `\dFp+` is used, a full description of each parser is shown, including the underlying functions and the list of recognized token types.

`\dFt[+] [pattern]`

Lists text search templates. If *pattern* is specified, only templates whose names match the pattern are shown. If the form `\dFt+` is used, additional information is shown about each template, including the underlying function names.

`\dg[S+] [pattern]`

Lists database roles. (Since the concepts of “users” and “groups” have been unified into “roles”, this command is now equivalent to `\du`.) By default, only user-created roles are shown; supply the *S* modifier to include system roles. If *pattern* is specified, only those roles whose names match the pattern are listed. If the form `\dg+` is used, additional information is shown about each role; currently this adds the comment for each role.

`\dl[+]`

This is an alias for `\lo_list`, which shows a list of large objects. If *+* is appended to the command name, each large object is listed with its associated permissions, if any.

`\dL[S+] [pattern]`

Lists procedural languages. If *pattern* is specified, only languages whose names match the pattern are listed. By default, only user-created languages are shown; supply the *S* modifier to include system objects. If *+* is appended to the command name, each language is listed with its call handler, validator, access privileges, and whether it is a system object.

`\dn[S+] [pattern]`

Lists schemas (namespaces). If *pattern* is specified, only schemas whose names match the pattern are listed. By default, only user-created objects are shown; supply a pattern or the *S* modifier to include system objects. If *+* is appended to the command name, each object is listed with its associated permissions, description, and the security officer, if any.


```
\do[S+] [ pattern [ arg_pattern [ arg_pattern ] ] ]
```

Lists operators with their operand and result types. If *pattern* is specified, only operators whose names match the pattern are listed. If one *arg_pattern* is specified, only prefix operators whose right argument's type name matches that pattern are listed. If two *arg_patterns* are specified, only binary operators whose argument type names match those patterns are listed. (Alternatively, write `-` for the unused argument of a unary operator.) By default, only user-created objects are shown; supply a pattern or the `s` modifier to include system objects. If `+` is appended to the command name, additional information about each operator is shown, currently just the name of the underlying function.

```
\do[S+] [ pattern ]
```

Lists collations. If *pattern* is specified, only collations whose names match the pattern are listed. By default, only user-created objects are shown; supply a pattern or the `s` modifier to include system objects. If `+` is appended to the command name, each collation is listed with its associated description, if any. Note that only collations usable with the current database's encoding are shown, so the results may vary in different databases of the same installation.

```
\dp[S] [ pattern ]
```

Lists tables, views and sequences with their associated access privileges. If *pattern* is specified, only tables, views and sequences whose names match the pattern are listed. By default only user-created objects are shown; supply a pattern or the `s` modifier to include system objects.

The [GRANT](#) and [REVOKE](#) commands are used to set access privileges. The meaning of the privilege display is explained in [Section 5.8](#).

```
\dP[itn+] [ pattern ]
```

Lists partitioned relations. If *pattern* is specified, only entries whose name matches the pattern are listed. The modifiers `t` (tables) and `i` (indexes) can be appended to the command, filtering the kind of relations to list. By default, partitioned tables and indexes are listed.

If the modifier `n` ("nested") is used, or a pattern is specified, then non-root partitioned relations are included, and a column is shown displaying the parent of each partitioned relation.

If `+` is appended to the command name, the sum of the sizes of each relation's partitions is also displayed, along with the relation's description. If `n` is combined with `+`, two sizes are shown: one including the total size of directly-attached leaf partitions, and another showing the total size of all partitions, including indirectly attached sub-partitions.

```
\drds [ role-pattern [ database-pattern ] ]
```

Lists defined configuration settings. These settings can be role-specific, database-specific, or both. *role-pattern* and *database-pattern* are used to select specific roles and databases to list, respectively. If omitted, or if `*` is specified, all settings are listed, including those not role-specific or database-specific, respectively.

The [ALTER ROLE](#) and [ALTER DATABASE](#) commands are used to define per-role and per-database configuration settings.

```
\drg[S] [ pattern ]
```

Lists information about each granted role membership, including assigned options (`ADMIN`, `INHERIT` and/or `SET`) and grantor. See the [GRANT](#) command for information about role memberships.

By default, only grants to user-created roles are shown; supply the `s` modifier to include system roles. If *pattern* is specified, only grants to those roles whose names match the pattern are listed.

```
\dRp[+] [ pattern ]
```

Lists replication publications. If *pattern* is specified, only those publications whose names match the pattern are listed. If `+` is appended to the command name, the tables and schemas associated with each publication are shown as well.

`\dRs[+] [pattern]`

Lists replication subscriptions. If *pattern* is specified, only those subscriptions whose names match the pattern are listed. If + is appended to the command name, additional properties of the subscriptions are shown.

`\duP [pattern]`

Lists profiles. If *pattern* is specified, only those profiles whose names match the pattern are listed. This command is only available for superusers.

`\dT[S+] [pattern]`

Lists data types. If *pattern* is specified, only types whose names match the pattern are listed. If + is appended to the command name, each type is listed with its internal name and size, its allowed values if it is an `enum` type, and its associated permissions. By default, only user-created objects are shown; supply a pattern or the *s* modifier to include system objects.

`\du[S+] [pattern]`

Lists database roles. (Since the concepts of “users” and “groups” have been unified into “roles”, this command is now equivalent to `\dg`.) By default, only user-created roles are shown; supply the *s* modifier to include system roles. If *pattern* is specified, only those roles whose names match the pattern are listed. If the form `\du+` is used, additional information is shown about each role; currently this adds the comment for each role.

`\dx[+] [pattern]`

Lists installed extensions. If *pattern* is specified, only those extensions whose names match the pattern are listed. If the form `\dx+` is used, all the objects belonging to each matching extension are listed.

`\dX [pattern]`

Lists extended statistics. If *pattern* is specified, only those extended statistics whose names match the pattern are listed.

The status of each kind of extended statistics is shown in a column named after its statistic kind (e.g. `Ndistinct`). `defined` means that it was requested when creating the statistics, and `NULL` means it wasn't requested. You can use `pg_stats_ext` if you'd like to know whether `ANALYZE` was run and statistics are available to the planner.

`\dy[+] [pattern]`

Lists event triggers. If *pattern* is specified, only those event triggers whose names match the pattern are listed. If + is appended to the command name, each object is listed with its associated description.

`\e or \edit [filename] [line_number]`

If *filename* is specified, the file is edited; after the editor exits, the file's content is copied into the current query buffer. If no *filename* is given, the current query buffer is copied to a temporary file which is then edited in the same fashion. Or, if the current query buffer is empty, the most recently executed query is copied to a temporary file and edited in the same fashion.

If you edit a file or the previous query, and you quit the editor without modifying the file, the query buffer is cleared. Otherwise, the new contents of the query buffer are re-parsed according to the normal rules of psql, treating the whole buffer as a single line. Any complete queries are immediately executed; that is, if the query buffer contains or ends with a semicolon, everything up to that point is executed and removed from the query buffer. Whatever remains in the query buffer is redisplayed. Type semicolon or `\g` to send it, or `\r` to cancel it by clearing the query buffer.

Treating the buffer as a single line primarily affects meta-commands: whatever is in the buffer after a meta-command will be taken as argument(s) to the meta-command, even if it spans multiple lines. (Thus you cannot make meta-command-using scripts this way. Use `\i` for that.)

If a line number is specified, psql will position the cursor on the specified line of the file or query buffer. Note that if a single all-digits argument is given, psql assumes it is a line number, not a file name.

Tip

See [Environment](#), below, for how to configure and customize your editor.

```
\echo text [ ... ]
```

Prints the evaluated arguments to standard output, separated by spaces and followed by a newline. This can be useful to intersperse information in the output of scripts. For example:

```
=> \echo `date`  
Tue Oct 26 21:40:57 CEST 1999
```

If the first argument is an unquoted `-n` the trailing newline is not written (nor is the first argument).

Tip

If you use the `\o` command to redirect your query output you might wish to use `\qecho` instead of this command. See also `\warn`.

```
\ef [ function_description [ line_number ] ]
```

This command fetches and edits the definition of the named function or procedure, in the form of a `CREATE OR REPLACE FUNCTION` or `CREATE OR REPLACE PROCEDURE` command. Editing is done in the same way as for `\edit`. If you quit the editor without saving, the statement is discarded. If you save and exit the editor, the updated command is executed immediately if you added a semicolon to it. Otherwise it is redisplayed; type semicolon or `\g` to send it, or `\r` to cancel.

The target function can be specified by name alone, or by name and arguments, for example `foo(integer, text)`. The argument types must be given if there is more than one function of the same name.

If no function is specified, a blank `CREATE FUNCTION` template is presented for editing.

If a line number is specified, psql will position the cursor on the specified line of the function body. (Note that the function body typically does not begin on the first line of the file.)

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\ef`, and neither variable interpolation nor backquote expansion are performed in the arguments.

Tip

See [Environment](#), below, for how to configure and customize your editor.

```
\encoding [ encoding ]
```

Sets the client character set encoding. Without an argument, this command shows the current encoding.

```
\errverbose
```

Repeats the most recent server error message at maximum verbosity, as though `VERBOSITY` were set to `verbose` and `SHOW_CONTEXT` were set to `always`.

```
\ev [ view_name [ line_number ] ]
```

This command fetches and edits the definition of the named view, in the form of a `CREATE OR REPLACE VIEW` command. Editing is done in the same way as for `\edit`. If you quit the editor without saving, the statement is discarded. If you save and exit the editor, the updated command is executed immediately if you added a semicolon to it. Otherwise it is redisplayed; type semicolon or `\g` to send it, or `\r` to cancel.

If no view is specified, a blank `CREATE VIEW` template is presented for editing.

If a line number is specified, psql will position the cursor on the specified line of the view definition.

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\ev`, and neither variable interpolation nor backquote expansion are performed in the arguments.

```
\f [ string ]
```

Sets the field separator for unaligned query output. The default is the vertical bar (`|`). It is equivalent to `\pset fieldsep`.

```
\g [ (option=value [...]) ] [ filename ]  
\g [ (option=value [...]) ] [ |command ]
```

Sends the current query buffer to the server for execution.

If parentheses appear after `\g`, they surround a space-separated list of `option=value` formatting-option clauses, which are interpreted in the same way as `\pset option value` commands, but take effect only for the duration of this query. In this list, spaces are not allowed around `=` signs, but are required between option clauses. If `=value` is omitted, the named `option` is changed in the same way as for `\pset option` with no explicit `value`.

If a `filename` or `|command` argument is given, the query's output is written to the named file or piped to the given shell command, instead of displaying it as usual. The file or command is written to only if the query successfully returns zero or more tuples, not if the query fails or is a non-data-returning SQL command.

If the current query buffer is empty, the most recently sent query is re-executed instead. Except for that behavior, `\g` without any arguments is essentially equivalent to a semicolon. With arguments, `\g` provides a “one-shot” alternative to the `\o` command, and additionally allows one-shot adjustments of the output formatting options normally set by `\pset`.

When the last argument begins with `|`, the entire remainder of the line is taken to be the `command` to execute, and neither variable interpolation nor backquote expansion are performed in it. The rest of the line is simply passed literally to the shell.

```
\gdesc
```

Shows the description (that is, the column names and data types) of the result of the current query buffer. The query is not actually executed; however, if it contains some type of syntax error, that error will be reported in the normal way.

If the current query buffer is empty, the most recently sent query is described instead.

```
\getenv psql_var env_var
```

Gets the value of the environment variable `env_var` and assigns it to the psql variable `psql_var`. If `env_var` is not defined in the psql process's environment, `psql_var` is not changed. Example:

```
=> \getenv home HOME  
=> \echo :home  
/home/postgres
```

\gexec

Sends the current query buffer to the server, then treats each column of each row of the query's output (if any) as an SQL statement to be executed. For example, to create an index on each column of `my_table`:

```
=> SELECT format('create index on my_table(%I)', attname)
-> FROM pg_attribute
-> WHERE attrelid = 'my_table'::regclass AND attnum > 0
-> ORDER BY attnum
-> \gexec
CREATE INDEX
CREATE INDEX
CREATE INDEX
CREATE INDEX
```

The generated queries are executed in the order in which the rows are returned, and left-to-right within each row if there is more than one column. NULL fields are ignored. The generated queries are sent literally to the server for processing, so they cannot be psql meta-commands nor contain psql variable references. If any individual query fails, execution of the remaining queries continues unless `ON_ERROR_STOP` is set. Execution of each query is subject to `ECHO` processing. (Setting `ECHO` to all or queries is often advisable when using `\gexec`.) Query logging, single-step mode, timing, and other query execution features apply to each generated query as well.

If the current query buffer is empty, the most recently sent query is re-executed instead.

\gset [*prefix*]

Sends the current query buffer to the server and stores the query's output into psql variables (see [Variables](#) below). The query to be executed must return exactly one row. Each column of the row is stored into a separate variable, named the same as the column. For example:

```
=> SELECT 'hello' AS var1, 10 AS var2
-> \gset
=> \echo :var1 :var2
hello 10
```

If you specify a *prefix*, that string is prepended to the query's column names to create the variable names to use:

```
=> SELECT 'hello' AS var1, 10 AS var2
-> \gset result_
=> \echo :result_var1 :result_var2
hello 10
```

If a column result is NULL, the corresponding variable is unset rather than being set.

If the query fails or does not return one row, no variables are changed.

If the current query buffer is empty, the most recently sent query is re-executed instead.

\gx [(*option=value* [...])] [*filename*]
\gx [(*option=value* [...])] [|*command*]

`\gx` is equivalent to `\g`, except that it forces expanded output mode for this query, as if `expanded=on` were included in the list of `\pset` options. See also `\x`.

\h or **\help** [*command*]

Gives syntax help on the specified SQL command. If *command* is not specified, then psql will list all the commands for which syntax help is available. If *command* is an asterisk (*), then syntax help on all SQL commands is shown.

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\help`, and neither variable interpolation nor backquote expansion are performed in the arguments.

Note

To simplify typing, commands that consists of several words do not have to be quoted. Thus it is fine to type `\help alter table`.

`\H` or `\html`

Turns on HTML query output format. If the HTML format is already on, it is switched back to the default aligned text format. This command is for compatibility and convenience, but see `\pset` about setting other output options.

`\i` or `\include filename [name=value] [value] [...]`

Reads input from the file *filename* and executes it as though it had been typed on the keyboard.

If *filename* is `-` (hyphen), then standard input is read until an EOF indication or `\q` meta-command. This can be used to intersperse interactive input with input from files. Note that Readline behavior will be used only if it is active at the outermost level.

When the optional arguments are passed, psql parses them for variable substitution. If *name=value* is passed, a named variable is created. If only the *value* is specified, it is saved as a positional variable named `PSQL_ARG1`, `PSQL_ARG2`, `PSQL_ARG3`, etc. These parameter variables are temporary and will be deleted after the script execution. Example:

```
\i expfile.sql filename=myfile.dmp lines=12
```

You can mix named and positional parameters, for example:

```
\i expfile.sql filename=myfile.dmp lines=12 1 100
```

In the above example, four variables will be available in the `expfile.sql` script:

- `filename` with the value `myfile.dmp`
- `lines` with the value `12`
- `PSQL_ARG1` with the value `1`
- `PSQL_ARG2` with the value `100`

If there are other psql variables with the same name, the temporary variable overrides the existing one for the time of script execution. If a script calls another one, both of them have their own sets of temporary variables, and the caller's variables cannot be used by another script (though they are not deleted yet).

Note

If you want to see the lines on the screen as they are read you must set the variable `ECHO` to `all`.

```
\if expression
\elif expression
\else
\endif
```

This group of commands implements nestable conditional blocks. A conditional block must begin with an `\if` and end with an `\endif`. In between there may be any number of `\elif` clauses, which

may optionally be followed by a single `\else` clause. Ordinary queries and other types of backslash commands may (and usually do) appear between the commands forming a conditional block.

The `\if` and `\elif` commands read their argument(s) and evaluate them as a Boolean expression. If the expression yields `true` then processing continues normally; otherwise, lines are skipped until a matching `\elif`, `\else`, or `\endif` is reached. Once an `\if` or `\elif` test has succeeded, the arguments of later `\elif` commands in the same block are not evaluated but are treated as false. Lines following an `\else` are processed only if no earlier matching `\if` or `\elif` succeeded.

The *expression* argument of an `\if` or `\elif` command is subject to variable interpolation and backquote expansion, just like any other backslash command argument. After that it is evaluated like the value of an on/off option variable. So a valid value is any unambiguous case-insensitive match for one of: `true`, `false`, `1`, `0`, `on`, `off`, `yes`, `no`. For example, `t`, `T`, and `tR` will all be considered to be `true`.

Expressions that do not properly evaluate to true or false will generate a warning and be treated as false.

Lines being skipped are parsed normally to identify queries and backslash commands, but queries are not sent to the server, and backslash commands other than conditionals (`\if`, `\elif`, `\else`, `\endif`) are ignored. Conditional commands are checked only for valid nesting. Variable references in skipped lines are not expanded, and backquote expansion is not performed either.

All the backslash commands of a given conditional block must appear in the same source file. If EOF is reached on the main input file or an `\include`-ed file before all local `\if`-blocks have been closed, then psql will raise an error.

Here is an example:

```
-- check for the existence of two separate records in the database and store
-- the results in separate psql variables
SELECT
    EXISTS(SELECT 1 FROM customer WHERE customer_id = 123) as is_customer,
    EXISTS(SELECT 1 FROM employee WHERE employee_id = 456) as is_employee
\gset
\if :is_customer
    SELECT * FROM customer WHERE customer_id = 123;
\elif :is_employee
    \echo 'is not a customer but is an employee'
    SELECT * FROM employee WHERE employee_id = 456;
\else
    \if yes
        \echo 'not a customer or employee'
    \else
        \echo 'this will never print'
    \endif
\endif
\ir or \include_relative filename
```

The `\ir` command is similar to `\i`, but resolves relative file names differently. When executing in interactive mode, the two commands behave identically. However, when invoked from a script, `\ir` interprets file names relative to the directory in which the script is located, rather than the current working directory.

```
\l[+] or \list[+] [ pattern ]
```

List the databases in the server and show their names, owners, character set encodings, and access privileges. If *pattern* is specified, only databases whose names match the pattern are listed. If `+` is appended to the command name, database sizes, default tablespaces, and descriptions are also displayed. (Size information is only available for databases that the current user can connect to.)

```
\lo_export loid filename
```

Reads the large object with OID *loid* from the database and writes it to *filename*. Note that this is subtly different from the server function `lo_export`, which acts with the permissions of the user that the database server runs as and on the server's file system.

Tip

Use `\lo_list` to find out the large object's OID.

```
\lo_import filename [ comment ]
```

Stores the file into a Postgres Pro large object. Optionally, it associates the given comment with the object. Example:

```
foo=> \lo_import '/home/peter/pictures/photo.xcf' 'a picture of me'
lo_import 152801
```

The response indicates that the large object received object ID 152801, which can be used to access the newly-created large object in the future. For the sake of readability, it is recommended to always associate a human-readable comment with every object. Both OIDs and comments can be viewed with the `\lo_list` command.

Note that this command is subtly different from the server-side `lo_import` because it acts as the local user on the local file system, rather than the server's user and file system.

```
\lo_list[+]
```

Shows a list of all Postgres Pro large objects currently stored in the database, along with any comments provided for them. If `+` is appended to the command name, each large object is listed with its associated permissions, if any.

```
\lo_unlink loid
```

Deletes the large object with OID *loid* from the database.

Tip

Use `\lo_list` to find out the large object's OID.

```
\o or \out [ filename ]
```

```
\o or \out [ |command ]
```

Arranges to save future query results to the file *filename* or pipe future results to the shell command *command*. If no argument is specified, the query output is reset to the standard output.

If the argument begins with `|`, then the entire remainder of the line is taken to be the *command* to execute, and neither variable interpolation nor backquote expansion are performed in it. The rest of the line is simply passed literally to the shell.

“Query results” includes all tables, command responses, and notices obtained from the database server, as well as output of various backslash commands that query the database (such as `\d`); but not error messages.

Tip

To intersperse text output in between query results, use `\qecho`.

`\p` or `\print`

Print the current query buffer to the standard output. If the current query buffer is empty, the most recently executed query is printed instead.

`\password [username]`

Changes the password of the specified user (by default, the current user). This command prompts for the new password, encrypts it, and sends it to the server as an `ALTER ROLE` command. This makes sure that the new password does not appear in cleartext in the command history, the server log, or elsewhere.

`\prompt [text] name`

Prompts the user to supply text, which is assigned to the variable *name*. An optional prompt string, *text*, can be specified. (For multiword prompts, surround the text with single quotes.)

By default, `\prompt` uses the terminal for input and output. However, if the `-f` command line switch was used, `\prompt` uses standard input and standard output.

`\pset [option [value]]`

This command sets options affecting the output of query result tables. *option* indicates which option is to be set. The semantics of *value* vary depending on the selected option. For some options, omitting *value* causes the option to be toggled or unset, as described under the particular option. If no such behavior is mentioned, then omitting *value* just results in the current setting being displayed.

`\pset` without any arguments displays the current status of all printing options.

Adjustable printing options are:

`border`

The *value* must be a number. In general, the higher the number the more borders and lines the tables will have, but details depend on the particular format. In HTML format, this will translate directly into the `border=...` attribute. In most other formats only values 0 (no border), 1 (internal dividing lines), and 2 (table frame) make sense, and values above 2 will be treated the same as `border = 2`. The `latex` and `latex-longtable` formats additionally allow a value of 3 to add dividing lines between data rows.

`columns`

Sets the target width for the `wrapped` format, and also the width limit for determining whether output is wide enough to require the pager or switch to the vertical display in expanded auto mode. Zero (the default) causes the target width to be controlled by the environment variable `COLUMNS`, or the detected screen width if `COLUMNS` is not set. In addition, if `columns` is zero then the `wrapped` format only affects screen output. If `columns` is nonzero then file and pipe output is wrapped to that width as well.

`csv_fieldsep`

Specifies the field separator to be used in CSV output format. If the separator character appears in a field's value, that field is output within double quotes, following standard CSV rules. The default is a comma.

`expanded` (or `x`)

If *value* is specified it must be either `on` or `off`, which will enable or disable expanded mode, or `auto`. If *value* is omitted the command toggles between the on and off settings. When expanded mode is enabled, query results are displayed in two columns, with the column name on the left and the data on the right. This mode is useful if the data wouldn't fit on the screen in the normal "horizontal" mode. In the auto setting, the expanded mode is used whenever the query output

has more than one column and is wider than the screen; otherwise, the regular mode is used. The auto setting is only effective in the aligned and wrapped formats. In other formats, it always behaves as if the expanded mode is off.

fieldsep

Specifies the field separator to be used in unaligned output format. That way one can create, for example, tab-separated output, which other programs might prefer. To set a tab as field separator, type `\pset fieldsep '\t'`. The default field separator is `'|'` (a vertical bar).

fieldsep_zero

Sets the field separator to use in unaligned output format to a zero byte.

footer

If *value* is specified it must be either `on` or `off` which will enable or disable display of the table footer (the `(n rows)` count). If *value* is omitted the command toggles footer display on or off.

format

Sets the output format to one of `aligned`, `asciidoc`, `csv`, `html`, `latex`, `latex-longtable`, `troff-ms`, `unaligned`, or `wrapped`. Unique abbreviations are allowed.

`aligned` format is the standard, human-readable, nicely formatted text output; this is the default.

`unaligned` format writes all columns of a row on one line, separated by the currently active field separator. This is useful for creating output that might be intended to be read in by other programs, for example, tab-separated or comma-separated format. However, the field separator character is not treated specially if it appears in a column's value; so CSV format may be better suited for such purposes.

`csv` format writes column values separated by commas, applying the quoting rules described in [RFC 4180](#). This output is compatible with the CSV format of the server's `COPY` command. A header line with column names is generated unless the `tuples_only` parameter is `on`. Titles and footers are not printed. Each row is terminated by the system-dependent end-of-line character, which is typically a single newline (`\n`) for Unix-like systems or a carriage return and newline sequence (`\r\n`) for Microsoft Windows. Field separator characters other than comma can be selected with `\pset csv_fieldsep`.

`wrapped` format is like `aligned` but wraps wide data values across lines to make the output fit in the target column width. The target width is determined as described under the `columns` option. Note that `psql` will not attempt to wrap column header titles; therefore, `wrapped` format behaves the same as `aligned` if the total width needed for column headers exceeds the target.

The `asciidoc`, `html`, `latex`, `latex-longtable`, and `troff-ms` formats put out tables that are intended to be included in documents using the respective mark-up language. They are not complete documents! This might not be necessary in HTML, but in LaTeX you must have a complete document wrapper. The `latex` format uses LaTeX's `tabular` environment. The `latex-longtable` format requires the LaTeX `longtable` and `booktabs` packages.

linestyle

Sets the border line drawing style to one of `ascii`, `old-ascii`, or `unicode`. Unique abbreviations are allowed. (That would mean one letter is enough.) The default setting is `ascii`. This option only affects the `aligned` and `wrapped` output formats.

`ascii` style uses plain ASCII characters. Newlines in data are shown using a `+` symbol in the right-hand margin. When the `wrapped` format wraps data from one line to the next without a newline character, a dot (`.`) is shown in the right-hand margin of the first line, and again in the left-hand margin of the following line.

`old-ascii` style uses plain ASCII characters, using the formatting style used in PostgreSQL 8.4 and earlier. Newlines in data are shown using a `:` symbol in place of the left-hand column separator. When the data is wrapped from one line to the next without a newline character, a `;` symbol is used in place of the left-hand column separator.

`unicode` style uses Unicode box-drawing characters. Newlines in data are shown using a carriage return symbol in the right-hand margin. When the data is wrapped from one line to the next without a newline character, an ellipsis symbol is shown in the right-hand margin of the first line, and again in the left-hand margin of the following line.

When the `border` setting is greater than zero, the `linestyle` option also determines the characters with which the border lines are drawn. Plain ASCII characters work everywhere, but Unicode characters look nicer on displays that recognize them.

`null`

Sets the string to be printed in place of a null value. The default is to print nothing, which can easily be mistaken for an empty string. For example, one might prefer `\pset null '(null)'`.

`numericlocale`

If *value* is specified it must be either `on` or `off` which will enable or disable display of a locale-specific character to separate groups of digits to the left of the decimal marker. If *value* is omitted the command toggles between regular and locale-specific numeric output.

`pager`

Controls use of a pager program for query and psql help output. When the `pager` option is `off`, the pager program is not used. When the `pager` option is `on`, the pager is used when appropriate, i.e., when the output is to a terminal and will not fit on the screen. The `pager` option can also be set to `always`, which causes the pager to be used for all terminal output regardless of whether it fits on the screen. `\pset pager` without a *value* toggles pager use on and off.

If the environment variable `PSQL_PAGER` or `PAGER` is set, output to be paged is piped to the specified program. Otherwise a platform-dependent default program (such as `more`) is used.

When using the `\watch` command to execute a query repeatedly, the environment variable `PSQL_WATCH_PAGER` is used to find the pager program instead, on Unix systems. This is configured separately because it may confuse traditional pagers, but can be used to send output to tools that understand psql's output format (such as `pspg --stream`).

`pager_min_lines`

If `pager_min_lines` is set to a number greater than the page height, the pager program will not be called unless there are at least this many lines of output to show. The default setting is 0.

`recordsep`

Specifies the record (line) separator to use in unaligned output format. The default is a newline character.

`recordsep_zero`

Sets the record separator to use in unaligned output format to a zero byte.

`tableattr` (or `T`)

In HTML format, this specifies attributes to be placed inside the `table` tag. This could for example be `cellpadding` or `bgcolor`. Note that you probably don't want to specify `border` here, as that is already taken care of by `\pset border`. If no *value* is given, the table attributes are unset.

In `latex-longtable` format, this controls the proportional width of each column containing a left-aligned data type. It is specified as a whitespace-separated list of values, e.g., `'0.2 0.2 0.6'`. Unspecified output columns use the last specified value.

title (or C)

Sets the table title for any subsequently printed tables. This can be used to give your output descriptive tags. If no *value* is given, the title is unset.

tuples_only (or t)

If *value* is specified it must be either `on` or `off` which will enable or disable tuples-only mode. If *value* is omitted the command toggles between regular and tuples-only output. Regular output includes extra information such as column headers, titles, and various footers. In tuples-only mode, only actual table data is shown.

unicode_border_linestyle

Sets the border drawing style for the `unicode` line style to one of `single` or `double`.

unicode_column_linestyle

Sets the column drawing style for the `unicode` line style to one of `single` or `double`.

unicode_header_linestyle

Sets the header drawing style for the `unicode` line style to one of `single` or `double`.

xheader_width

Sets the maximum width of the header for expanded output to one of `full` (the default value), `column`, `page`, or an *integer value*.

`full`: the expanded header is not truncated, and will be as wide as the widest output line.

`column`: truncate the header line to the width of the first column.

`page`: truncate the header line to the terminal width.

integer value: specify the exact maximum width of the header line.

Illustrations of how these different formats look can be seen in [Examples](#), below.

Tip

There are various shortcut commands for `\pset`. See `\a`, `\C`, `\f`, `\H`, `\t`, `\T`, and `\x`.

\q or \quit

Quits the psql program. In a script file, only execution of that script is terminated.

\qecho text [...]

This command is identical to `\echo` except that the output will be written to the query output channel, as set by `\o`.

\r or \reset

Resets (clears) the query buffer.

\s [filename]

Print psql's command line history to *filename*. If *filename* is omitted, the history is written to the standard output (using the pager if appropriate). This command is not available if psql was built without Readline support.

```
\set [ name [ value [ ... ] ] ]
```

Sets the psql variable *name* to *value*, or if more than one value is given, to the concatenation of all of them. If only one argument is given, the variable is set to an empty-string value. To unset a variable, use the `\unset` command.

`\set` without any arguments displays the names and values of all currently-set psql variables.

Valid variable names can contain letters, digits, and underscores. See [Variables](#) below for details. Variable names are case-sensitive.

Certain variables are special, in that they control psql's behavior or are automatically set to reflect connection state. These variables are documented in [Variables](#), below.

Note

This command is unrelated to the SQL command `SET`.

```
\setenv name [ value ]
```

Sets the environment variable *name* to *value*, or if the *value* is not supplied, unsets the environment variable. Example:

```
testdb=> \setenv PAGER less
testdb=> \setenv LESS -imx4F
```

```
\sf[+] function_description
```

This command fetches and shows the definition of the named function or procedure, in the form of a `CREATE OR REPLACE FUNCTION` or `CREATE OR REPLACE PROCEDURE` command. The definition is printed to the current query output channel, as set by `\o`.

The target function can be specified by name alone, or by name and arguments, for example `foo(integer, text)`. The argument types must be given if there is more than one function of the same name.

If `+` is appended to the command name, then the output lines are numbered, with the first line of the function body being line 1.

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\sf`, and neither variable interpolation nor backquote expansion are performed in the arguments.

```
\sv[+] view_name
```

This command fetches and shows the definition of the named view, in the form of a `CREATE OR REPLACE VIEW` command. The definition is printed to the current query output channel, as set by `\o`.

If `+` is appended to the command name, then the output lines are numbered from 1.

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\sv`, and neither variable interpolation nor backquote expansion are performed in the arguments.

```
\t
```

Toggles the display of output column name headings and row count footer. This command is equivalent to `\pset tuples_only` and is provided for convenience.

```
\T table_options
```

Specifies attributes to be placed within the `table` tag in HTML output format. This command is equivalent to `\pset tableattr table_options`.

`\timing [on | off]`

With a parameter, turns displaying of how long each SQL statement takes on or off. Without a parameter, toggles the display between on and off. The display is in milliseconds; intervals longer than 1 second are also shown in minutes:seconds format, with hours and days fields added if needed.

`\unset name`

Unsets (deletes) the psql variable *name*.

Most variables that control psql's behavior cannot be unset; instead, an `\unset` command is interpreted as setting them to their default values. See [Variables](#) below.

`\w or \write filename`

`\w or \write |command`

Writes the current query buffer to the file *filename* or pipes it to the shell command *command*. If the current query buffer is empty, the most recently executed query is written instead.

If the argument begins with `|`, then the entire remainder of the line is taken to be the *command* to execute, and neither variable interpolation nor backquote expansion are performed in it. The rest of the line is simply passed literally to the shell.

`\warn text [...]`

This command is identical to `\echo` except that the output will be written to psql's standard error channel, rather than standard output.

`\watch [i[nterval]=seconds] [c[ount]=times] [m[in_rows]=rows] [seconds]`

Repeatedly execute the current query buffer (as `\g` does) until interrupted, or the query fails, or the execution count limit (if given) is reached, or the query no longer returns the minimum number of rows. Wait the specified number of seconds (default 2) between executions. For backwards compatibility, *seconds* can be specified with or without an *interval=* prefix. Each query result is displayed with a header that includes the `\pset title` string (if any), the time as of query start, and the delay interval.

If the current query buffer is empty, the most recently sent query is re-executed instead.

`\x [on | off | auto]`

Sets or toggles expanded table formatting mode. As such it is equivalent to `\pset expanded`.

`\z[S] [pattern]`

Lists tables, views and sequences with their associated access privileges. If a *pattern* is specified, only tables, views and sequences whose names match the pattern are listed. By default only user-created objects are shown; supply a pattern or the *S* modifier to include system objects.

This is an alias for `\dp` ("display privileges").

`\! [command]`

With no argument, escapes to a sub-shell; psql resumes when the sub-shell exits. With an argument, executes the shell command *command*.

Unlike most other meta-commands, the entire remainder of the line is always taken to be the argument(s) of `\!`, and neither variable interpolation nor backquote expansion are performed in the arguments. The rest of the line is simply passed literally to the shell.

`\? [topic]`

Shows help information. The optional *topic* parameter (defaulting to *commands*) selects which part of psql is explained: *commands* describes psql's backslash commands; *options* describes the com-

mand-line options that can be passed to psql; and `variables` shows help about psql configuration variables.

`\;`

Backslash-semicolon is not a meta-command in the same way as the preceding commands; rather, it simply causes a semicolon to be added to the query buffer without any further processing.

Normally, psql will dispatch an SQL command to the server as soon as it reaches the command-ending semicolon, even if more input remains on the current line. Thus for example entering

```
select 1; select 2; select 3;
```

will result in the three SQL commands being individually sent to the server, with each one's results being displayed before continuing to the next command. However, a semicolon entered as `\;` will not trigger command processing, so that the command before it and the one after are effectively combined and sent to the server in one request. So for example

```
select 1\; select 2\; select 3;
```

results in sending the three SQL commands to the server in a single request, when the non-backslashed semicolon is reached. The server executes such a request as a single transaction, unless there are explicit `BEGIN/COMMIT` commands included in the string to divide it into multiple transactions. (See [Section 56.2.2.1](#) for more details about how the server handles multi-query strings.)

Patterns

The various `\d` commands accept a *pattern* parameter to specify the object name(s) to be displayed. In the simplest case, a pattern is just the exact name of the object. The characters within a pattern are normally folded to lower case, just as in SQL names; for example, `\dt FOO` will display the table named `foo`. As in SQL names, placing double quotes around a pattern stops folding to lower case. Should you need to include an actual double quote character in a pattern, write it as a pair of double quotes within a double-quote sequence; again this is in accord with the rules for SQL quoted identifiers. For example, `\dt "FOO" "BAR"` will display the table named `FOO"BAR` (not `foo"bar`). Unlike the normal rules for SQL names, you can put double quotes around just part of a pattern, for instance `\dt FOO"FOO"BAR` will display the table named `fooFOObar`.

Whenever the *pattern* parameter is omitted completely, the `\d` commands display all objects that are visible in the current schema search path — this is equivalent to using `*` as the pattern. (An object is said to be *visible* if its containing schema is in the search path and no object of the same kind and name appears earlier in the search path. This is equivalent to the statement that the object can be referenced by name without explicit schema qualification.) To see all objects in the database regardless of visibility, use `*.*` as the pattern.

Within a pattern, `*` matches any sequence of characters (including no characters) and `?` matches any single character. (This notation is comparable to Unix shell file name patterns.) For example, `\dt int*` displays tables whose names begin with `int`. But within double quotes, `*` and `?` lose these special meanings and are just matched literally.

A relation pattern that contains a dot (`.`) is interpreted as a schema name pattern followed by an object name pattern. For example, `\dt foo*.*bar*` displays all tables whose table name includes `bar` that are in schemas whose schema name starts with `foo`. When no dot appears, then the pattern matches only objects that are visible in the current schema search path. Again, a dot within double quotes loses its special meaning and is matched literally. A relation pattern that contains two dots (`..`) is interpreted as a database name followed by a schema name pattern followed by an object name pattern. The database name portion will not be treated as a pattern and must match the name of the currently connected database, else an error will be raised.

A schema pattern that contains a dot (`.`) is interpreted as a database name followed by a schema name pattern. For example, `\dn mydb.*foo*` displays all schemas whose schema name includes `foo`. The data-

base name portion will not be treated as a pattern and must match the name of the currently connected database, else an error will be raised.

Advanced users can use regular-expression notations such as character classes, for example `[0-9]` to match any digit. All regular expression special characters work as specified in [Section 9.7.3](#), except for `.` which is taken as a separator as mentioned above, `*` which is translated to the regular-expression notation `.*`, `?` which is translated to `.`, and `$` which is matched literally. You can emulate these pattern characters at need by writing `?` for `.`, `(R+|)` for `R*`, or `(R|)` for `R?`. `$` is not needed as a regular-expression character since the pattern must match the whole name, unlike the usual interpretation of regular expressions (in other words, `$` is automatically appended to your pattern). Write `*` at the beginning and/or end if you don't wish the pattern to be anchored. Note that within double quotes, all regular expression special characters lose their special meanings and are matched literally. Also, the regular expression special characters are matched literally in operator name patterns (i.e., the argument of `\do`).

Advanced Features

Variables

psql provides variable substitution features similar to common Unix command shells. Variables are simply name/value pairs, where the value can be any string of any length. The name must consist of letters (including non-Latin letters), digits, and underscores.

To set a variable, use the psql meta-command `\set`. For example,

```
testdb=> \set foo bar
```

sets the variable `foo` to the value `bar`. To retrieve the content of the variable, precede the name with a colon, for example:

```
testdb=> \echo :foo
bar
```

This works in both regular SQL commands and meta-commands; there is more detail in [SQL Interpolation](#), below.

If you call `\set` without a second argument, the variable is set to an empty-string value. To unset (i.e., delete) a variable, use the command `\unset`. To show the values of all variables, call `\set` without any argument.

Note

The arguments of `\set` are subject to the same substitution rules as with other commands. Thus you can construct interesting references such as `\set :foo 'something'` and get “soft links” or “variable variables” of Perl or PHP fame, respectively. Unfortunately (or fortunately?), there is no way to do anything useful with these constructs. On the other hand, `\set bar :foo` is a perfectly valid way to copy a variable.

A number of these variables are treated specially by psql. They represent certain option settings that can be changed at run time by altering the value of the variable, or in some cases represent changeable state of psql. By convention, all specially treated variables' names consist of all upper-case ASCII letters (and possibly digits and underscores). To ensure maximum compatibility in the future, avoid using such variable names for your own purposes.

Variables that control psql's behavior generally cannot be unset or set to invalid values. An `\unset` command is allowed but is interpreted as setting the variable to its default value. A `\set` command without a second argument is interpreted as setting the variable to `on`, for control variables that accept that value, and is rejected for others. Also, control variables that accept the values `on` and `off` will also accept other common spellings of Boolean values, such as `true` and `false`.

The specially treated variables are:

AUTOCOMMIT

When `on` (the default), each SQL command is automatically committed upon successful completion. To postpone commit in this mode, you must enter a `BEGIN` or `START TRANSACTION` SQL command. When `off` or `unset`, SQL commands are not committed until you explicitly issue `COMMIT` or `END`. The autocommit-off mode works by issuing an implicit `BEGIN` for you, just before any command that is not already in a transaction block and is not itself a `BEGIN` or other transaction-control command, nor a command that cannot be executed inside a transaction block (such as `VACUUM`).

Note

In autocommit-off mode, you must explicitly abandon any failed transaction by entering `ABORT` or `ROLLBACK`. Also keep in mind that if you exit the session without committing, your work will be lost.

Note

The autocommit-on mode is Postgres Pro's traditional behavior, but autocommit-off is closer to the SQL spec. If you prefer autocommit-off, you might wish to set it in the system-wide `psqlrc` file or your `~/.psqlrc` file.

COMP_KEYWORD_CASE

Determines which letter case to use when completing an SQL key word. If set to `lower` or `upper`, the completed word will be in lower or upper case, respectively. If set to `preserve-lower` or `preserve-upper` (the default), the completed word will be in the case of the word already entered, but words being completed without anything entered will be in lower or upper case, respectively.

DBNAME

The name of the database you are currently connected to. This is set every time you connect to a database (including program start-up), but can be changed or unset.

ECHO

If set to `all`, all nonempty input lines are printed to standard output as they are read. (This does not apply to lines read interactively.) To select this behavior on program start-up, use the switch `-a`. If set to `queries`, `psql` prints each query to standard output as it is sent to the server. The switch to select this behavior is `-e`. If set to `errors`, then only failed queries are displayed on standard error output. The switch for this behavior is `-b`. If set to `none` (the default), then no queries are displayed.

ECHO_HIDDEN

When this variable is set to `on` and a backslash command queries the database, the query is first shown. This feature helps you to study Postgres Pro internals and provide similar functionality in your own programs. (To select this behavior on program start-up, use the switch `-E`.) If you set this variable to the value `noexec`, the queries are just shown but are not actually sent to the server and executed. The default value is `off`.

ENCODING

The current client character set encoding. This is set every time you connect to a database (including program start-up), and when you change the encoding with `\encoding`, but it can be changed or unset.

ERROR

`true` if the last SQL query failed, `false` if it succeeded. See also `SQLSTATE`.

FETCH_COUNT

If this variable is set to an integer value greater than zero, the results of `SELECT` queries are fetched and displayed in groups of that many rows, rather than the default behavior of collecting the entire result set before display. Therefore only a limited amount of memory is used, regardless of the size of the result set. Settings of 100 to 1000 are commonly used when enabling this feature. Keep in mind that when using this feature, a query might fail after having already displayed some rows.

Tip

Although you can use any output format with this feature, the default `aligned` format tends to look bad because each group of `FETCH_COUNT` rows will be formatted separately, leading to varying column widths across the row groups. The other output formats work better.

HIDE_TABLEAM

If this variable is set to `true`, a table's access method details are not displayed. This is mainly useful for regression tests.

HIDE_TOAST_COMPRESSION

If this variable is set to `true`, column compression method details are not displayed. This is mainly useful for regression tests.

HISTCONTROL

If this variable is set to `ignoreSpace`, lines which begin with a space are not entered into the history list. If set to a value of `ignoreDups`, lines matching the previous history line are not entered. A value of `ignoreBoth` combines the two options. If set to `none` (the default), all lines read in interactive mode are saved on the history list.

Note

This feature was shamelessly plagiarized from Bash.

HISTFILE

The file name that will be used to store the history list. If unset, the file name is taken from the `PSQL_HISTORY` environment variable. If that is not set either, the default is `~/.psql_history`, or `%APPDATA%\postgresql\psql_history` on Windows. For example, putting:

```
\set HISTFILE ~/.psql_history-:DBNAME
```

in `~/.psqlrc` will cause `psql` to maintain a separate history for each database.

Note

This feature was shamelessly plagiarized from Bash.

HISTSIZE

The maximum number of commands to store in the command history (default 500). If set to a negative value, no limit is applied.

Note

This feature was shamelessly plagiarized from Bash.

HOST

The database server host you are currently connected to. This is set every time you connect to a database (including program start-up), but can be changed or unset.

IGNOREEOF

If set to 1 or less, sending an EOF character (usually **Control+D**) to an interactive session of psql will terminate the application. If set to a larger numeric value, that many consecutive EOF characters must be typed to make an interactive session terminate. If the variable is set to a non-numeric value, it is interpreted as 10. The default is 0.

Note

This feature was shamelessly plagiarized from Bash.

LASTOID

The value of the last affected OID, as returned from an `INSERT` or `\lo_import` command. This variable is only guaranteed to be valid until after the result of the next SQL command has been displayed. Postgres Pro servers since version 12 do not support OID system columns anymore, thus `LASTOID` will always be 0 following `INSERT` when targeting such servers.

LAST_ERROR_MESSAGE**LAST_ERROR_SQLSTATE**

The primary error message and associated `SQLSTATE` code for the most recent failed query in the current psql session, or an empty string and 00000 if no error has occurred in the current session.

ON_ERROR_ROLLBACK

When set to `on`, if a statement in a transaction block generates an error, the error is ignored and the transaction continues. When set to `interactive`, such errors are only ignored in interactive sessions, and not when reading script files. When set to `off` (the default), a statement in a transaction block that generates an error aborts the entire transaction. The error rollback mode works by issuing an implicit `SAVEPOINT` for you, just before each command that is in a transaction block, and then rolling back to the savepoint if the command fails.

ON_ERROR_STOP

By default, command processing continues after an error. When this variable is set to `on`, processing will instead stop immediately. In interactive mode, psql will return to the command prompt; otherwise, psql will exit, returning error code 3 to distinguish this case from fatal error conditions, which are reported using error code 1. In either case, any currently running scripts (the top-level script, if any, and any other scripts which it may have in invoked) will be terminated immediately. If the top-level command string contained multiple SQL commands, processing will stop with the current command.

PORT

The database server port to which you are currently connected. This is set every time you connect to a database (including program start-up), but can be changed or unset.

PROMPT1**PROMPT2****PROMPT3**

These specify what the prompts psql issues should look like. See [Prompting](#) below.

QUIET

Setting this variable to `on` is equivalent to the command line option `-q`. It is probably not too useful in interactive mode.

ROW_COUNT

The number of rows returned or affected by the last SQL query, or 0 if the query failed or did not report a row count.

SERVER_VERSION_NAME**SERVER_VERSION_NUM**

The server's version number as a string, for example 9.6.2, 10.1 or 11beta1, and in numeric form, for example 90602 or 100001. These are set every time you connect to a database (including program start-up), but can be changed or unset.

SHELL_ERROR

true if the last shell command failed, false if it succeeded. This applies to shell commands invoked via the \!, \g, \o, \w, and \copy meta-commands, as well as backquote (`) expansion. Note that for \o, this variable is updated when the output pipe is closed by the next \o command. See also SHELL_EXIT_CODE.

SHELL_EXIT_CODE

The exit status returned by the last shell command. 0-127 represent program exit codes, 128-255 indicate termination by a signal, and -1 indicates failure to launch a program or to collect its exit status. This applies to shell commands invoked via the \!, \g, \o, \w, and \copy meta-commands, as well as backquote (`) expansion. Note that for \o, this variable is updated when the output pipe is closed by the next \o command. See also SHELL_ERROR.

SHOW_ALL_RESULTS

When this variable is set to off, only the last result of a combined query (\;) is shown instead of all of them. The default is on. The off behavior is for compatibility with older versions of psql.

SHOW_CONTEXT

This variable can be set to the values never, errors, or always to control whether CONTEXT fields are displayed in messages from the server. The default is errors (meaning that context will be shown in error messages, but not in notice or warning messages). This setting has no effect when VERBOSITY is set to terse or sqlstate. (See also \errverbose, for use when you want a verbose version of the error you just got.)

SINGLELINE

Setting this variable to on is equivalent to the command line option -S.

SINGLESTEP

Setting this variable to on is equivalent to the command line option -s.

SQLSTATE

The error code (see [Appendix A](#)) associated with the last SQL query's failure, or 00000 if it succeeded.

USER

The database user you are currently connected as. This is set every time you connect to a database (including program start-up), but can be changed or unset.

VERBOSITY

This variable can be set to the values default, verbose, terse, or sqlstate to control the verbosity of error reports. (See also \errverbose, for use when you want a verbose version of the error you just got.)

```
VERSION
VERSION_NAME
VERSION_NUM
```

These variables are set at program start-up to reflect psql's version, respectively as a verbose string, a short string (e.g., 9.6.2, 10.1, or 11beta1), and a number (e.g., 90602 or 100001). They can be changed or unset.

SQL Interpolation

A key feature of psql variables is that you can substitute (“interpolate”) them into regular SQL statements, as well as the arguments of meta-commands. Furthermore, psql provides facilities for ensuring that variable values used as SQL literals and identifiers are properly quoted. The syntax for interpolating a value without any quoting is to prepend the variable name with a colon (:). For example,

```
testdb=> \set foo 'my_table'
testdb=> SELECT * FROM :foo;
```

would query the table `my_table`. Note that this may be unsafe: the value of the variable is copied literally, so it can contain unbalanced quotes, or even backslash commands. You must make sure that it makes sense where you put it.

When a value is to be used as an SQL literal or identifier, it is safest to arrange for it to be quoted. To quote the value of a variable as an SQL literal, write a colon followed by the variable name in single quotes. To quote the value as an SQL identifier, write a colon followed by the variable name in double quotes. These constructs deal correctly with quotes and other special characters embedded within the variable value. The previous example would be more safely written this way:

```
testdb=> \set foo 'my_table'
testdb=> SELECT * FROM :"foo";
```

Variable interpolation will not be performed within quoted SQL literals and identifiers. Therefore, a construction such as `':foo'` doesn't work to produce a quoted literal from a variable's value (and it would be unsafe if it did work, since it wouldn't correctly handle quotes embedded in the value).

One example use of this mechanism is to copy the contents of a file into a table column. First load the file into a variable and then interpolate the variable's value as a quoted string:

```
testdb=> \set content `cat my_file.txt`
testdb=> INSERT INTO my_table VALUES (: 'content');
```

(Note that this still won't work if `my_file.txt` contains NUL bytes. psql does not support embedded NUL bytes in variable values.)

Since colons can legally appear in SQL commands, an apparent attempt at interpolation (that is, `:name`, `:'name'`, or `:"name"`) is not replaced unless the named variable is currently set. In any case, you can escape a colon with a backslash to protect it from substitution.

The `: {?name}` special syntax returns TRUE or FALSE depending on whether the variable exists or not, and is thus always substituted, unless the colon is backslash-escaped.

The colon syntax for variables is standard SQL for embedded query languages, such as ECPG. The colon syntaxes for array slices and type casts are Postgres Pro extensions, which can sometimes conflict with the standard usage. The colon-quote syntax for escaping a variable's value as an SQL literal or identifier is a psql extension.

Prompting

The prompts psql issues can be customized to your preference. The three variables `PROMPT1`, `PROMPT2`, and `PROMPT3` contain strings and special escape sequences that describe the appearance of the prompt. Prompt 1 is the normal prompt that is issued when psql requests a new command. Prompt 2 is issued when more input is expected during command entry, for example because the command was not terminated with a semicolon or a quote was not closed. Prompt 3 is issued when you are running an SQL `COPY FROM STDIN` command and you need to type in a row value on the terminal.

The value of the selected prompt variable is printed literally, except where a percent sign (%) is encountered. Depending on the next character, certain other text is substituted instead. Defined substitutions are:

%M

The full host name (with domain name) of the database server, or [local] if the connection is over a Unix domain socket, or [local:/dir/name], if the Unix domain socket is not at the compiled in default location.

%m

The host name of the database server, truncated at the first dot, or [local] if the connection is over a Unix domain socket.

%>

The port number at which the database server is listening.

%n

The database session user name. (The expansion of this value might change during a database session as the result of the command `SET SESSION AUTHORIZATION`.)

%/

The name of the current database.

%~

Like %/, but the output is ~ (tilde) if the database is your default database.

%#

If the session user is a database superuser, then a #, otherwise a >. (The expansion of this value might change during a database session as the result of the command `SET SESSION AUTHORIZATION`.)

%p

The process ID of the backend currently connected to.

%R

In prompt 1 normally =, but @ if the session is in an inactive branch of a conditional block, or ^ if in single-line mode, or ! if the session is disconnected from the database (which can happen if `\connect` fails). In prompt 2 %R is replaced by a character that depends on why psql expects more input: - if the command simply wasn't terminated yet, but * if there is an unfinished `/* ... */` comment, a single quote if there is an unfinished quoted string, a double quote if there is an unfinished quoted identifier, a dollar sign if there is an unfinished dollar-quoted string, or (if there is an unmatched left parenthesis. In prompt 3 %R doesn't produce anything.

%x

Transaction status: an empty string when not in a transaction block, or * when in a transaction block, or ! when in a failed transaction block, or ? when the transaction state is indeterminate (for example, because there is no connection).

%l

The line number inside the current statement, starting from 1.

%*digits*

The character with the indicated octal code is substituted.

`%:name:`

The value of the psql variable *name*. See [Variables](#), above, for details.

`%`command``

The output of *command*, similar to ordinary “back-tick” substitution.

`%[... %]`

Prompts can contain terminal control characters which, for example, change the color, background, or style of the prompt text, or change the title of the terminal window. In order for the line editing features of Readline to work properly, these non-printing control characters must be designated as invisible by surrounding them with `%[` and `%)`. Multiple pairs of these can occur within the prompt. For example:

```
testdb=> \set PROMPT1 '%[%033[1;33;40m%]n@%/%R%[%033[0m%]## '
```

results in a boldfaced (1;) yellow-on-black (33;40) prompt on VT100-compatible, color-capable terminals.

`%w`

Whitespace of the same width as the most recent output of `PROMPT1`. This can be used as a `PROMPT2` setting, so that multi-line statements are aligned with the first line, but there is no visible secondary prompt.

To insert a percent sign into your prompt, write `%%`. The default prompts are `'%/%R%x%# '` for prompts 1 and 2, and `'>> '` for prompt 3.

Note

This feature was shamelessly plagiarized from `tcsh`.

Command-Line Editing

psql uses the Readline or libedit library, if available, for convenient line editing and retrieval. The command history is automatically saved when psql exits and is reloaded when psql starts up. Type up-arrow or control-P to retrieve previous lines.

You can also use tab completion to fill in partially-typed keywords and SQL object names in many (by no means all) contexts. For example, at the start of a command, typing `ins` and pressing TAB will fill in `insert into`. Then, typing a few characters of a table or schema name and pressing TAB will fill in the unfinished name, or offer a menu of possible completions when there's more than one. (Depending on the library in use, you may need to press TAB more than once to get a menu.)

Tab completion for SQL object names requires sending queries to the server to find possible matches. In some contexts this can interfere with other operations. For example, after `BEGIN` it will be too late to issue `SET TRANSACTION ISOLATION LEVEL` if a tab-completion query is issued in between. If you do not want tab completion at all, you can turn it off permanently by putting this in a file named `.inputrc` in your home directory:

```
$if psql
set disable-completion on
$endif
```

(This is not a psql but a Readline feature. Read its documentation for further details.)

The `-n` (`--no-readline`) command line option can also be useful to disable use of Readline for a single run of psql. This prevents tab completion, use or recording of command line history, and editing of multi-line commands. It is particularly useful when you need to copy-and-paste text that contains TAB characters.

Environment

COLUMNS

If `\pset columns` is zero, controls the width for the `wrapped` format and width for determining if wide output requires the pager or should be switched to the vertical format in expanded auto mode.

PGDATABASE

PGHOST

PGPORT

PGUSER

Default connection parameters (see [Section 35.15](#)).

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

PSQL_EDITOR

EDITOR

VISUAL

Editor used by the `\e`, `\ef`, and `\ev` commands. These variables are examined in the order listed; the first that is set is used. If none of them is set, the default is to use `vi` on Unix systems or `notepad.exe` on Windows systems.

PSQL_EDITOR_LINENUMBER_ARG

When `\e`, `\ef`, or `\ev` is used with a line number argument, this variable specifies the command-line argument used to pass the starting line number to the user's editor. For editors such as Emacs or `vi`, this is a plus sign. Include a trailing space in the value of the variable if there needs to be space between the option name and the line number. Examples:

```
PSQL_EDITOR_LINENUMBER_ARG='+'
```

```
PSQL_EDITOR_LINENUMBER_ARG='--line '
```

The default is `+` on Unix systems (corresponding to the default editor `vi`, and useful for many other common editors); but there is no default on Windows systems.

PSQL_HISTORY

Alternative location for the command history file. Tilde (`~`) expansion is performed.

PSQL_PAGER

PAGER

If a query's results do not fit on the screen, they are piped through this command. Typical values are `more` or `less`. Use of the pager can be disabled by setting `PSQL_PAGER` or `PAGER` to an empty string, or by adjusting the pager-related options of the `\pset` command. These variables are examined in the order listed; the first that is set is used. If neither of them is set, the default is to use `more` on most platforms, but `less` on Cygwin.

PSQL_WATCH_PAGER

When a query is executed repeatedly with the `\watch` command, a pager is not used by default. This behavior can be changed by setting `PSQL_WATCH_PAGER` to a pager command, on Unix systems. The `pspg` pager (not part of PostgreSQL but available in many open source software distributions) can display the output of `\watch` if started with the option `--stream`.

PSQLRC

Alternative location of the user's `.psqlrc` file. Tilde (`~`) expansion is performed.

SHELL

Command executed by the `\!` command.

TMPDIR

Directory for storing temporary files. The default is `/tmp`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Files

psqlrc and `~/.psqlrc`

Unless it is passed an `-x` option, `psql` attempts to read and execute commands from the system-wide startup file (`psqlrc`) and then the user's personal startup file (`~/.psqlrc`), after connecting to the database but before accepting normal commands. These files can be used to set up the client and/or the server to taste, typically with `\set` and `SET` commands.

The system-wide startup file is named `psqlrc`. By default it is sought in the installation's "system configuration" directory, which is most reliably identified by running `pg_config --sysconfdir`. Typically this directory will be `./etc/` relative to the directory containing the Postgres Pro executables. The directory to look in can be set explicitly via the `PGSYSCONFDIR` environment variable.

The user's personal startup file is named `.psqlrc` and is sought in the invoking user's home directory. On Windows the personal startup file is instead named `%APPDATA%\postgresql\psqlrc.conf`. In either case, this default file path can be overridden by setting the `PSQLRC` environment variable.

Both the system-wide startup file and the user's personal startup file can be made `psql`-version-specific by appending a dash and the Postgres Pro major or minor release identifier to the file name, for example `~/.psqlrc-17` or `~/.psqlrc-17.5.1`. The most specific version-matching file will be read in preference to a non-version-specific file. These version suffixes are added after determining the file path as explained above.

`.psql_history`

The command-line history is stored in the file `~/.psql_history`, or `%APPDATA%\postgresql\psql_history` on Windows.

The location of the history file can be set explicitly via the `HISTFILE` `psql` variable or the `PSQL_HISTORY` environment variable.

Notes

- `psql` works best with servers of the same or an older major version. Backslash commands are particularly likely to fail if the server is of a newer version than `psql` itself. However, backslash commands of the `\d` family should work with servers of versions back to 9.2, though not necessarily with servers newer than `psql` itself. The general functionality of running SQL commands and displaying query results should also work with servers of a newer major version, but this cannot be guaranteed in all cases.

If you want to use `psql` to connect to several servers of different major versions, it is recommended that you use the newest version of `psql`. Alternatively, you can keep around a copy of `psql` from each major version and be sure to use the version that matches the respective server. But in practice, this additional complication should not be necessary.

- Before Postgres Pro 9.6, the `-c` option implied `-x` (`--no-psqlrc`); this is no longer the case.
- Before PostgreSQL 8.4, `psql` allowed the first argument of a single-letter backslash command to start directly after the command, without intervening whitespace. Now, some whitespace is required.

Notes for Windows Users

psql is built as a “console application”. Since the Windows console windows use a different encoding than the rest of the system, you must take special care when using 8-bit characters within psql. If psql detects a problematic console code page, it will warn you at startup. To change the console code page, two things are necessary:

- Set the code page by entering `cmd.exe /c chcp 1252`. (1252 is a code page that is appropriate for German; replace it with your value.) If you are using Cygwin, you can put this command in `/etc/profile`.
- Set the console font to `Lucida Console`, because the raster font does not work with the ANSI code page.

By default, psql works in the UTF-8 encoding and uses Windows Unicode API for console output. The console code page must be set to 65001 to ensure the correct display of all characters supported by your Windows console font.

Postgres Pro Windows installer ships the `less.exe` pager with UTF-8 support and provides a shortcut that opens the console window with `Lucida Console` font and `codepage 65001` settings. If you use a different pager, make sure it also supports UTF-8.

You can override the default psql encoding by setting the `PGCLIENTENCODING` environment variable.

Examples

The first example shows how to spread a command over several lines of input. Notice the changing prompt:

```
testdb=> CREATE TABLE my_table (  
testdb(>   first integer not null default 0,  
testdb(>   second text)  
testdb-> ;  
CREATE TABLE
```

Now look at the table definition again:

```
testdb=> \d my_table  
          Table "public.my_table"  
Column | Type   | Collation | Nullable | Default  
-----+-----+-----+-----+-----  
first  | integer |           | not null | 0  
second | text    |           |          |
```

Now we change the prompt to something more interesting:

```
testdb=> \set PROMPT1 '%n%m %~%R%# '  
peter@localhost testdb=>
```

Let's assume you have filled the table with data and want to take a look at it:

```
peter@localhost testdb=> SELECT * FROM my_table;  
 first | second  
-----+-----  
    1  | one  
    2  | two  
    3  | three  
    4  | four  
(4 rows)
```

You can display tables in different ways by using the `\pset` command:

```
peter@localhost testdb=> \pset border 2  
Border style is 2.
```

```
peter@localhost testdb=> SELECT * FROM my_table;
```

```
+-----+-----+
| first | second |
+-----+-----+
|      1 | one     |
|      2 | two     |
|      3 | three   |
|      4 | four    |
+-----+-----+
(4 rows)
```

```
peter@localhost testdb=> \pset border 0
```

```
Border style is 0.
```

```
peter@localhost testdb=> SELECT * FROM my_table;
```

```
first second
-----
1 one
2 two
3 three
4 four
(4 rows)
```

```
peter@localhost testdb=> \pset border 1
```

```
Border style is 1.
```

```
peter@localhost testdb=> \pset format csv
```

```
Output format is csv.
```

```
peter@localhost testdb=> \pset tuples_only
```

```
Tuples only is on.
```

```
peter@localhost testdb=> SELECT second, first FROM my_table;
```

```
one,1
two,2
three,3
four,4
```

```
peter@localhost testdb=> \pset format unaligned
```

```
Output format is unaligned.
```

```
peter@localhost testdb=> \pset fieldsep '\t'
```

```
Field separator is "\t".
```

```
peter@localhost testdb=> SELECT second, first FROM my_table;
```

```
one      1
two      2
three    3
four     4
```

Alternatively, use the short commands:

```
peter@localhost testdb=> \a \t \x
```

```
Output format is aligned.
```

```
Tuples only is off.
```

```
Expanded display is on.
```

```
peter@localhost testdb=> SELECT * FROM my_table;
```

```
-[ RECORD 1 ]-
first  | 1
second | one
-[ RECORD 2 ]-
first  | 2
second | two
-[ RECORD 3 ]-
first  | 3
second | three
```

```
-[ RECORD 4 ]-
first  | 4
second | four
```

Also, these output format options can be set for just one query by using `\g`:

```
peter@localhost testdb=> SELECT * FROM my_table
peter@localhost testdb-> \g (format=aligned tuples_only=off expanded=on)
-[ RECORD 1 ]-
first  | 1
second | one
-[ RECORD 2 ]-
first  | 2
second | two
-[ RECORD 3 ]-
first  | 3
second | three
-[ RECORD 4 ]-
first  | 4
second | four
```

Here is an example of using the `\df` command to find only functions with names matching `int*pl` and whose second argument is of type `bigint`:

```
testdb=> \df int*pl * bigint
                                List of functions
 Schema | Name      | Result data type | Argument data types | Type
-----+-----+-----+-----+-----
pg_catalog | int28pl | bigint          | smallint, bigint   | func
pg_catalog | int48pl | bigint          | integer, bigint    | func
pg_catalog | int8pl  | bigint          | bigint, bigint     | func
(3 rows)
```

When suitable, query results can be shown in a crosstab representation with the `\crosstabview` command:

```
testdb=> SELECT first, second, first > 2 AS gt2 FROM my_table;
 first | second | gt2
-----+-----+-----
      1 | one    | f
      2 | two    | f
      3 | three  | t
      4 | four   | t
(4 rows)
```

```
testdb=> \crosstabview first second
 first | one | two | three | four
-----+-----+-----+-----+-----
      1 | f   |     |       |
      2 |     | f   |       |
      3 |     |     | t     |
      4 |     |     |       | t
(4 rows)
```

This second example shows a multiplication table with rows sorted in reverse numerical order and columns with an independent, ascending numerical order.

```
testdb=> SELECT t1.first as "A", t2.first+100 AS "B", t1.first*(t2.first+100) as "AxB",
testdb-> row_number() over(order by t2.first) AS ord
testdb-> FROM my_table t1 CROSS JOIN my_table t2 ORDER BY 1 DESC
testdb-> \crosstabview "A" "B" "AxB" ord
```

A	101	102	103	104
4	404	408	412	416
3	303	306	309	312
2	202	204	206	208
1	101	102	103	104

(4 rows)

reindexdb

reindexdb — reindex a Postgres Pro database

Synopsis

```
reindexdb [connection-option...] [option...] [ -S | --schema schema ] ... [ -t | --table table ] ... [ -i  
| --index index ] ... [ -s | --system ] [ dbname | -a | --all ]
```

Description

reindexdb is a utility for rebuilding indexes in a Postgres Pro database.

reindexdb is a wrapper around the SQL command [REINDEX](#). There is no effective difference between reindexing databases via this utility and via other methods for accessing the server.

Options

reindexdb accepts the following command-line arguments:

-a
--all

Reindex all databases.

--concurrently

Use the `CONCURRENTLY` option. See [REINDEX](#), where all the caveats of this option are explained in detail.

[-d] *dbname*
[--dbname=] *dbname*

Specifies the name of the database to be reindexed, when `-a/--all` is not used. If this is not specified, the database name is read from the environment variable `PGDATABASE`. If that is not set, the user name specified for the connection is used. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

-e
--echo

Echo the commands that reindexdb generates and sends to the server.

-i *index*
--index=*index*

Recreate *index* only. Multiple indexes can be recreated by writing multiple `-i` switches.

-j *njobs*
--jobs=*njobs*

Execute the reindex commands in parallel by running *njobs* commands simultaneously. This option may reduce the processing time but it also increases the load on the database server.

reindexdb will open *njobs* connections to the database, so make sure your [max_connections](#) setting is high enough to accommodate all connections.

Note that this option is incompatible with the `--system` option.

-q
--quiet

Do not display progress messages.

`-s`
`--system`

Reindex database's system catalogs only.

`-S schema`
`--schema=schema`

Reindex *schema* only. Multiple schemas can be reindexed by writing multiple `-S` switches.

`-t table`
`--table=table`

Reindex *table* only. Multiple tables can be reindexed by writing multiple `-t` switches.

`--tablespace=tablespace`

Specifies the tablespace where indexes are rebuilt. (This name is processed as a double-quoted identifier.)

`-v`
`--verbose`

Print detailed information during processing.

`-V`
`--version`

Print the reindexdb version and exit.

`-?`
`--help`

Show help about reindexdb command line arguments, and exit.

reindexdb also accepts the following command-line arguments for connection parameters:

`-h host`
`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`
`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U username`
`--username=username`

User name to connect as.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force reindexdb to prompt for a password before connecting to a database.

This option is never essential, since `reindexdb` will automatically prompt for a password if the server demands password authentication. However, `reindexdb` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

`--maintenance-db=dbname`

When the `-a/--all` is used, connect to this database to gather the list of databases to reindex. If not specified, the `postgres` database will be used, or if that does not exist, `template1` will be used. This can be a [connection string](#). If so, connection string parameters will override any conflicting command line options. Also, connection string parameters other than the database name itself will be re-used when connecting to other databases.

Environment

PGDATABASE
PGHOST
PGPORT
PGUSER

Default connection parameters

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [REINDEX](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To reindex the database `test`:

```
$ reindexdb test
```

To reindex the table `foo` and the index `bar` in a database named `abcd`:

```
$ reindexdb --table=foo --index=bar abcd
```

See Also

[REINDEX](#)

vacuumdb

vacuumdb — garbage-collect and analyze a Postgres Pro database

Synopsis

```
vacuumdb [connection-option...] [option...] [ -t | --table table [( column [,...] )] ] ... [ dbname | -a | --all ]
```

```
vacuumdb [connection-option...] [option...] [ -n | --schema schema ] ... [ dbname | -a | --all ]
```

```
vacuumdb [connection-option...] [option...] [ -N | --exclude-schema schema ] ... [ dbname | -a | --all ]
```

Description

vacuumdb is a utility for cleaning a Postgres Pro database. vacuumdb will also generate internal statistics used by the Postgres Pro query optimizer.

vacuumdb is a wrapper around the SQL command [VACUUM](#). There is no effective difference between vacuuming and analyzing databases via this utility and via other methods for accessing the server.

Options

vacuumdb accepts the following command-line arguments:

-a
--all

Vacuum all databases.

--buffer-usage-limit *size*

Specifies the [Buffer Access Strategy](#) ring buffer size for a given invocation of vacuumdb. This size is used to calculate the number of shared buffers which will be reused as part of this strategy. See [VACUUM](#).

[-d] *dbname*
[--dbname=] *dbname*

Specifies the name of the database to be cleaned or analyzed, when -a/--all is not used. If this is not specified, the database name is read from the environment variable PGDATABASE. If that is not set, the user name specified for the connection is used. The *dbname* can be a [connection string](#). If so, connection string parameters will override any conflicting command line options.

--disable-page-skipping

Disable skipping pages based on the contents of the visibility map.

-e
--echo

Echo the commands that vacuumdb generates and sends to the server.

-f
--full

Perform “full” vacuuming.

-F
--freeze

Aggressively “freeze” tuples.

`--force-index-cleanup`

Always remove index entries pointing to dead tuples.

`-j njobs`

`--jobs=njobs`

Execute the vacuum or analyze commands in parallel by running *njobs* commands simultaneously. This option may reduce the processing time but it also increases the load on the database server.

`vacuumdb` will open *njobs* connections to the database, so make sure your [max_connections](#) setting is high enough to accommodate all connections.

Note that using this mode together with the `-f` (FULL) option might cause deadlock failures if certain system catalogs are processed in parallel.

`--min-mxid-age mxid_age`

Only execute the vacuum or analyze commands on tables with a multixact ID age of at least *mxid_age*. This setting is useful for prioritizing tables to process to prevent multixact ID wraparound (see [Section 24.1.5.1](#)).

For the purposes of this option, the multixact ID age of a relation is the greatest of the ages of the main relation and its associated TOAST table, if one exists. Since the commands issued by `vacuumdb` will also process the TOAST table for the relation if necessary, it does not need to be considered separately.

`--min-xid-age xid_age`

Only execute the vacuum or analyze commands on tables with a transaction ID age of at least *xid_age*. This setting is useful for prioritizing tables to process to prevent transaction ID wraparound (see [Section 24.1.5](#)).

For the purposes of this option, the transaction ID age of a relation is the greatest of the ages of the main relation and its associated TOAST table, if one exists. Since the commands issued by `vacuumdb` will also process the TOAST table for the relation if necessary, it does not need to be considered separately.

`-n schema`

`--schema=schema`

Clean or analyze all tables in *schema* only. Multiple schemas can be vacuumed by writing multiple `-n` switches.

`-N schema`

`--exclude-schema=schema`

Do not clean or analyze any tables in *schema*. Multiple schemas can be excluded by writing multiple `-N` switches.

`--no-index-cleanup`

Do not remove index entries pointing to dead tuples.

`--no-process-main`

Skip the main relation.

`--no-process-toast`

Skip the TOAST table associated with the table to vacuum, if any.

`--no-truncate`

Do not truncate empty pages at the end of the table.

`-P parallel_workers`
`--parallel=parallel_workers`

Specify the number of parallel workers for *parallel vacuum*. This allows the vacuum to leverage multiple CPUs to process indexes. See [VACUUM](#).

`-q`
`--quiet`

Do not display progress messages.

`--skip-locked`

Skip relations that cannot be immediately locked for processing.

`-t table [(column [,...])]`
`--table=table [(column [,...])]`

Clean or analyze *table* only. Column names can be specified only in conjunction with the `--analyze` or `--analyze-only` options. Multiple tables can be vacuumed by writing multiple `-t` switches.

Tip

If you specify columns, you probably have to escape the parentheses from the shell. (See examples below.)

`-v`
`--verbose`

Print detailed information during processing.

`-V`
`--version`

Print the vacuumdb version and exit.

`-z`
`--analyze`

Also calculate statistics for use by the optimizer.

`-Z`
`--analyze-only`

Only calculate statistics for use by the optimizer (no vacuum).

`--analyze-in-stages`

Only calculate statistics for use by the optimizer (no vacuum), like `--analyze-only`. Run three stages of analyze; the first stage uses the lowest possible statistics target (see [default_statistics_target](#)) to produce usable statistics faster, and subsequent stages build the full statistics.

This option is only useful to analyze a database that currently has no statistics or has wholly incorrect ones, such as if it is newly populated from a restored dump or by `pg_upgrade`. Be aware that running with this option in a database with existing statistics may cause the query optimizer choices to become transiently worse due to the low statistics targets of the early stages.

`-?`
`--help`

Show help about vacuumdb command line arguments, and exit.

vacuumdb also accepts the following command-line arguments for connection parameters:

`-h host`
`--host=host`

Specifies the host name of the machine on which the server is running. If the value begins with a slash, it is used as the directory for the Unix domain socket.

`-p port`
`--port=port`

Specifies the TCP port or local Unix domain socket file extension on which the server is listening for connections.

`-U username`
`--username=username`

User name to connect as.

`-w`
`--no-password`

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

`-W`
`--password`

Force vacuumdb to prompt for a password before connecting to a database.

This option is never essential, since vacuumdb will automatically prompt for a password if the server demands password authentication. However, vacuumdb will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-W` to avoid the extra connection attempt.

`--maintenance-db=dbname`

When the `-a/--all` is used, connect to this database to gather the list of databases to vacuum. If not specified, the `postgres` database will be used, or if that does not exist, `template1` will be used. This can be a [connection string](#). If so, connection string parameters will override any conflicting command line options. Also, connection string parameters other than the database name itself will be re-used when connecting to other databases.

Environment

PGDATABASE
PGHOST
PGPORT
PGUSER

Default connection parameters

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

Diagnostics

In case of difficulty, see [VACUUM](#) and [psql](#) for discussions of potential problems and error messages. The database server must be running at the targeted host. Also, any default connection settings and environment variables used by the libpq front-end library will apply.

Examples

To clean the database `test`:

```
$ vacuumdb test
```

To clean and analyze for the optimizer a database named `bigdb`:

```
$ vacuumdb --analyze bigdb
```

To clean a single table `foo` in a database named `xyzyy`, and analyze a single column `bar` of the table for the optimizer:

```
$ vacuumdb --analyze --verbose --table='foo(bar)' xyzyy
```

To clean all tables in the `foo` and `bar` schemas in a database named `xyzyy`:

```
$ vacuumdb --schema='foo' --schema='bar' xyzyy
```

See Also

[VACUUM](#)

Postgres Pro Server Applications

This part contains reference information for Postgres Pro server applications and support utilities. These commands can only be run usefully on the host where the database server resides. Other utility programs are listed in [Postgres Pro Client Applications](#).

initdb

initdb — create a new Postgres Pro database cluster

Synopsis

```
initdb [option...] [ --pgdata | -D ] directory
```

Description

initdb creates a new Postgres Pro *database cluster*.

Creating a database cluster consists of creating the *directories* in which the cluster data will live, generating the shared catalog tables (tables that belong to the whole cluster rather than to any particular database), and creating the `postgres`, `template1`, and `template0` databases. The `postgres` database is a default database meant for use by users, utilities and third party applications. `template1` and `template0` are meant as source databases to be copied by later `CREATE DATABASE` commands. `template0` should never be modified, but you can add objects to `template1`, which by default will be copied into databases created later. See [Section 22.4](#) for more details.

Although `initdb` will attempt to create the specified data directory, it might not have permission if the parent directory of the desired data directory is root-owned. To initialize in such a setup, create an empty data directory as root, then use `chown` to assign ownership of that directory to the database user account, then `su` to become the database user to run `initdb`.

`initdb` must be run as the user that will own the server process, because the server needs to have access to the files and directories that `initdb` creates. Since the server cannot be run as root, you must not run `initdb` as root either. (It will in fact refuse to do so.)

For security reasons the new cluster created by `initdb` will only be accessible by the cluster owner by default. The `--allow-group-access` option allows any user in the same group as the cluster owner to read files in the cluster. This is useful for performing backups as a non-privileged user.

`initdb` initializes the database cluster's default locale and character set encoding. These can also be set separately for each database when it is created. `initdb` determines those settings for the template databases, which will serve as the default for all other databases.

By default, `initdb` takes the locale settings from the environment, and determines the encoding from the locale settings. If the locale provider is not specified, it is set automatically based on the `lc_collate` value for the cluster: `libc` is used for the C or POSIX locales, and `icu` is used for other locales.

To choose a different locale for the cluster, use the option `--locale`. There are also individual options `--lc-*` and `--icu-locale` (see below) to set values for the individual locale categories. Note that inconsistent settings for different locale categories can give nonsensical results, so this should be used with care.

Alternatively, `initdb` can use the ICU library to provide locale services by specifying `--locale-provider=icu`. The server must be built with ICU support. To choose the specific ICU locale ID to apply, use the option `--icu-locale`. Note that for implementation reasons and to support legacy code, `initdb` will still select and initialize `libc` locale settings when the ICU locale provider is used.

When `initdb` runs, it will print out the locale settings it has chosen. If you have complex requirements or specified multiple options, it is advisable to check that the result matches what was intended.

More details about locale settings can be found in [Section 23.1](#).

To alter the default encoding, use the `--encoding`. More details can be found in [Section 23.3](#).

Options

`-A authmethod`

`--auth=authmethod`

This option specifies the default authentication method for local users used in `pg_hba.conf` (host and local lines). See [Section 20.1](#) for an overview of valid values.

`initdb` will prepopulate `pg_hba.conf` entries using the specified authentication method for non-replication as well as replication connections.

Do not use `trust` unless you trust all local users on your system. `trust` is the default for ease of installation.

`--auth-host=authmethod`

This option specifies the authentication method for local users via TCP/IP connections used in `pg_hba.conf` (host lines).

`--auth-local=authmethod`

This option specifies the authentication method for local users via Unix-domain socket connections used in `pg_hba.conf` (local lines).

`-D directory`

`--pgdata=directory`

This option specifies the directory where the database cluster should be stored. This is the only information required by `initdb`, but you can avoid writing it by setting the `PGDATA` environment variable, which can be convenient since the database server (`postgres`) can find the data directory later by the same variable.

`-E encoding`

`--encoding=encoding`

Selects the encoding of the template databases. This will also be the default encoding of any database you create later, unless you override it then. The character sets supported by the Postgres Pro server are described in [Section 23.3.1](#).

By default, the template database encoding is derived from the locale. If `--no-locale` is specified (or equivalently, if the locale is `C` or `POSIX`), then the default is `UTF8` for the ICU provider and `SQL_ASCII` for the `libc` provider.

`-g`

`--allow-group-access`

Allows users in the same group as the cluster owner to read all cluster files created by `initdb`. This option is ignored on Windows as it does not support POSIX-style group permissions.

`--icu-locale=locale`

Specifies the ICU locale when the ICU provider is used. Locale support is described in [Section 23.1](#).

`--icu-rules=rules`

Specifies additional collation rules to customize the behavior of the default collation. This is supported for ICU only.

`-k`

`--data-checksums`

Use checksums on data pages to help detect corruption by the I/O system that would otherwise be silent. Enabling checksums may incur a noticeable performance penalty. If set, checksums are calculated for all objects, in all databases. All checksum failures will be reported in the `pg_stat_data-base` view. See [Section 28.2](#) for details.

By default, Postgres Pro clusters are initialized with checksums enabled. To change this behavior, provide the `--no-data-checksums` option.

`--no-data-checksums`

Disable checksums on data pages.

By default, Postgres Pro clusters are initialized with checksums enabled.

`--locale=locale`

Sets the default locale for the database cluster. If this option is not specified, the locale is inherited from the environment that `initdb` runs in. Locale support is described in [Section 23.1](#).

If `--locale-provider` is builtin, `--locale` or `--builtin-locale` must be specified and set to `C` or `C.UTF-8`.

`--lc-collate=locale`

`--lc-ctype=locale`

`--lc-messages=locale`

`--lc-monetary=locale`

`--lc-numeric=locale`

`--lc-time=locale`

Like `--locale`, but only sets the locale in the specified category. The `--lc-collate` option overrides the `--locale` setting, regardless of whether the collation provider is specified.

`--no-locale`

Equivalent to `--locale=C`.

`--builtin-locale=locale`

Specifies the locale name when the builtin provider is used. Locale support is described in [Section 23.1](#).

`--locale-provider={builtin|libc|icu}`

This option sets the locale provider for databases created in the new cluster. It can be overridden in the `CREATE DATABASE` command when new databases are subsequently created. See [Section 23.1.4](#).

`--pwfile=filename`

Makes `initdb` read the bootstrap superuser's password from a file. The first line of the file is taken as the password.

`--no-tune`

By default, the use of `pgpro_tune` is enabled in `initdb`. Use this option to disable it.

`--tune=options`

Sets parameters for `pgpro_tune`. See [their description](#) for further information.

`-T config`

`--text-search-config=config`

Sets the default text search configuration. See [default_text_search_config](#) for further information.

`-U username`

`--username=username`

Sets the user name of the *bootstrap superuser*. This defaults to the name of the operating-system user running `initdb`.

-W
--pwprompt

Makes `initdb` prompt for a password to give the bootstrap superuser. If you don't plan on using password authentication, this is not important. Otherwise you won't be able to use password authentication until you have a password set up.

-X *directory*
--waldir=*directory*

This option specifies the directory where the write-ahead log should be stored.

--wal-segsize=*size*

Set the *WAL segment size*, in megabytes. This is the size of each individual file in the WAL log. The default size is 16 megabytes. The value must be a power of 2 between 1 and 1024 (megabytes). This option can only be set during initialization, and cannot be changed later.

It may be useful to adjust this size to control the granularity of WAL log shipping or archiving. Also, in databases with a high volume of WAL, the sheer number of WAL files per directory can become a performance and management problem. Increasing the WAL file size will reduce the number of WAL files.

Other, less commonly used, options are also available:

-c *name=value*
--set *name=value*

Forcibly set the server parameter *name* to *value* during `initdb`, and also install that setting in the generated `postgresql.conf` file, so that it will apply during future server runs. This option can be given more than once to set several parameters. It is primarily useful when the environment is such that the server will not start at all using the default parameters.

-d
--debug

Print debugging output from the bootstrap backend and a few other messages of lesser interest for the general public. The bootstrap backend is the program `initdb` uses to create the catalog tables. This option generates a tremendous amount of extremely boring output.

--discard-caches

Run the bootstrap backend with the `debug_discard_caches=1` option. This takes a very long time and is only of use for deep debugging.

-L *directory*

Specifies where `initdb` should find its input files to initialize the database cluster. This is normally not necessary. You will be told if you need to specify their location explicitly.

-n
--no-clean

By default, when `initdb` determines that an error prevented it from completely creating the database cluster, it removes any files it might have created before discovering that it cannot finish the job. This option inhibits tidying-up and is thus useful for debugging.

-N
--no-sync

By default, `initdb` will wait for all files to be written safely to disk. This option causes `initdb` to return without waiting, which is faster, but means that a subsequent operating system crash can

leave the data directory corrupt. Generally, this option is useful for testing, but should not be used when creating a production installation.

`--no-instructions`

By default, `initdb` will write instructions for how to start the cluster at the end of its output. This option causes those instructions to be left out. This is primarily intended for use by tools that wrap `initdb` in platform-specific behavior, where those instructions are likely to be incorrect.

`-s`

`--show`

Show internal settings and exit, without doing anything else. This can be used to debug the `initdb` installation.

`--sync-method=method`

When set to `fsync`, which is the default, `initdb` will recursively open and synchronize all files in the data directory. The search for files will follow symbolic links for the WAL directory and each configured tablespace.

On Linux, `syncfs` may be used instead to ask the operating system to synchronize the whole file systems that contain the data directory, the WAL files, and each tablespace. See [recovery_init_sync_method](#) for information about the caveats to be aware of when using `syncfs`.

This option has no effect when `--no-sync` is used.

`-S`

`--sync-only`

Safely write all database files to disk and exit. This does not perform any of the normal `initdb` operations. Generally, this option is useful for ensuring reliable recovery after changing `fsync` from `off` to `on`.

Other options for testing 64-bit XIDs:

`-m START_MX_ID`

`--multixact-id=START_MX_ID`

Specifies the initial multixact ID value in the decimal format for a new database instance to test database upgrades, default value is 0.

`-o START_MX_OFFSET`

`--multixact-offset=START_MX_OFFSET`

Specifies the initial multixact offset value in the decimal format for a new database instance to test database upgrades, default value is 0.

`-x START_XID`

`--xid=START_XID`

Specifies the initial XID value in the decimal format for a new database instance to test database upgrades, default value is 0.

Other options:

`-V`

`--version`

Print the `initdb` version and exit.

`-?`

`--help`

Show help about `initdb` command line arguments, and exit.

Environment

PGDATA

Specifies the directory where the database cluster is to be stored; can be overridden using the `-D` option.

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

PGPRO_TUNE

Specifies whether to use [pgpro_tune](#) without modifying command-line options directly. If set to `disable`, this environment variable has the same effect as passing `--no-tune`, disabling the use of `pgpro_tune`. If set to any other value, and neither `--tune` nor `--no-tune` is explicitly specified on the command line, the behavior is the same as if `--tune` is specified with the value of `PGPRO_TUNE`.

TZ

Specifies the default time zone of the created database cluster. The value should be a full time zone name (see [Section 8.5.3](#)).

Notes

`initdb` can also be invoked via `pg_ctl initdb`.

See Also

[pg_ctl](#), [postgres](#), [Section 20.1](#)

pg_archivecleanup

pg_archivecleanup — clean up Postgres Pro WAL archive files

Synopsis

```
pg_archivecleanup [option...] archivelocation oldestkeptwalfile
```

Description

pg_archivecleanup is designed to be used as an `archive_cleanup_command` to clean up WAL file archives when running as a standby server (see [Section 26.2](#)). pg_archivecleanup can also be used as a standalone program to clean WAL file archives.

To configure a standby server to use pg_archivecleanup, put this into its `postgresql.conf` configuration file:

```
archive_cleanup_command = 'pg_archivecleanup archivelocation %r'
```

where *archivelocation* is the directory from which WAL segment files should be removed.

When used within [archive_cleanup_command](#), all WAL files logically preceding the value of the `%r` argument will be removed from *archivelocation*. This minimizes the number of files that need to be retained, while preserving crash-restart capability. Use of this parameter is appropriate if the *archivelocation* is a transient staging area for this particular standby server, but *not* when the *archivelocation* is intended as a long-term WAL archive area, or when multiple standby servers are recovering from the same archive location.

When used as a standalone program all WAL files logically preceding the *oldestkeptwalfile* will be removed from *archivelocation*. In this mode, if you specify a `.partial` or `.backup` file name, then only the file prefix will be used as the *oldestkeptwalfile*. This treatment of `.backup` file name allows you to remove all WAL files archived prior to a specific base backup without error. For example, the following example will remove all files older than WAL file name `00000001000000037000000010`:

```
pg_archivecleanup -d archive 00000001000000037000000010.00000020.backup
```

```
pg_archivecleanup: keep WAL file "archive/00000001000000037000000010" and later
pg_archivecleanup: removing file "archive/0000000100000003700000000F"
pg_archivecleanup: removing file "archive/0000000100000003700000000E"
```

pg_archivecleanup assumes that *archivelocation* is a directory readable and writable by the server-owning user.

Options

pg_archivecleanup accepts the following command-line arguments:

`-b`

`--clean-backup-history`

Remove backup history files as well.

`-d`

`--debug`

Print lots of debug logging output on `stderr`.

`-n`

`--dry-run`

Print the names of the files that would have been removed on `stdout` (performs a dry run).

`-V`
`--version`

Print the `pg_archivecleanup` version and exit.

`-x extension`
`--strip-extension=extension`

Provide an extension that will be stripped from all file names before deciding if they should be deleted. This is typically useful for cleaning up archives that have been compressed during storage, and therefore have had an extension added by the compression program. For example: `-x .gz`.

`-?`
`--help`

Show help about `pg_archivecleanup` command line arguments, and exit.

Environment

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

`pg_archivecleanup` is designed to work with PostgreSQL 8.0 and later when used as a standalone utility, or with PostgreSQL 9.0 and later when used as an archive cleanup command.

`pg_archivecleanup` is written in C and has an easy-to-modify source code, with specifically designated sections to modify for your own needs

Examples

On Linux or Unix systems, you might use:

```
archive_cleanup_command = 'pg_archivecleanup -d /mnt/standby/archive %r 2>>cleanup.log'
```

where the archive directory is physically located on the standby server, so that the `archive_command` is accessing it across NFS, but the files are local to the standby. This will:

- produce debugging output in `cleanup.log`
- remove no-longer-needed files from the archive directory

pg_checksums

pg_checksums — enable, disable or check data checksums in a Postgres Pro database cluster

Synopsis

```
pg_checksums [option...] [[ -D | --pgdata ]datadir]
```

Description

pg_checksums checks, enables or disables data checksums in a Postgres Pro cluster. The server must be shut down cleanly before running pg_checksums. When verifying checksums, the exit status is zero if there are no checksum errors, and nonzero if at least one checksum failure is detected. When enabling or disabling checksums, the exit status is nonzero if the operation failed.

When verifying checksums, every file in the cluster is scanned. When enabling checksums, each relation file block with a changed checksum is rewritten in-place. Disabling checksums only updates the file pg_control.

Options

The following command-line options are available:

`-D directory`
`--pgdata=directory`

Specifies the directory where the database cluster is stored.

`-c`
`--check`

Checks checksums. This is the default mode if nothing else is specified.

`-d`
`--disable`

Disables checksums.

`-e`
`--enable`

Enables checksums.

`-f filenode`
`--filenode=filenode`

Only validate checksums in the relation with filenode *filenode*.

`-N`
`--no-sync`

By default, pg_checksums will wait for all files to be written safely to disk. This option causes pg_checksums to return without waiting, which is faster, but means that a subsequent operating system crash can leave the updated data directory corrupt. Generally, this option is useful for testing but should not be used on a production installation. This option has no effect when using `--check`.

`-P`
`--progress`

Enable progress reporting. Turning this on will deliver a progress report while checking or enabling checksums.

`--sync-method=method`

When set to `fsync`, which is the default, `pg_checksums` will recursively open and synchronize all files in the data directory. The search for files will follow symbolic links for the WAL directory and each configured tablespace.

On Linux, `syncfs` may be used instead to ask the operating system to synchronize the whole file systems that contain the data directory, the WAL files, and each tablespace. See [recovery_init_sync_method](#) for information about the caveats to be aware of when using `syncfs`.

This option has no effect when `--no-sync` is used.

`-v`

`--verbose`

Enable verbose output. Lists all checked files.

`-V`

`--version`

Print the `pg_checksums` version and exit.

`-?`

`--help`

Show help about `pg_checksums` command line arguments, and exit.

Environment

`PGDATA`

Specifies the directory where the database cluster is stored; can be overridden using the `-D` option.

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

Enabling checksums in a large cluster can potentially take a long time. During this operation, the cluster or other programs that write to the data directory must not be started or else data loss may occur.

When using a replication setup with tools which perform direct copies of relation file blocks (for example [pg_rewind](#)), enabling or disabling checksums can lead to page corruptions in the shape of incorrect checksums if the operation is not done consistently across all nodes. When enabling or disabling checksums in a replication setup, it is thus recommended to stop all the clusters before switching them all consistently. Destroying all standbys, performing the operation on the primary and finally recreating the standbys from scratch is also safe.

If `pg_checksums` is aborted or killed while enabling or disabling checksums, the cluster's data checksum configuration remains unchanged, and `pg_checksums` can be re-run to perform the same operation.

pg_createsubscriber

pg_createsubscriber — convert a physical replica into a new logical replica

Synopsis

```
pg_createsubscriber [option...] { -d | --database }dbname { -D | --pgdata }datadir { -P | --publisher-server }connstr
```

Description

pg_createsubscriber creates a new logical replica from a physical standby server. All tables in the specified database are included in the [logical replication](#) setup. A pair of publication and subscription objects are created for each database. It must be run at the target server.

After a successful run, the state of the target server is analogous to a fresh logical replication setup. The main difference between the logical replication setup and pg_createsubscriber is how the data synchronization is done. pg_createsubscriber does not copy the initial table data. It does only the synchronization phase, which ensures each table is brought up to a synchronized state.

pg_createsubscriber targets large database systems because in logical replication setup, most of the time is spent doing the initial data copy. Furthermore, a side effect of this long time spent synchronizing data is usually a large amount of changes to be applied (that were produced during the initial data copy), which increases even more the time when the logical replica will be available. For smaller databases, it is recommended to set up logical replication with initial data synchronization. For details, see the `CREATE SUBSCRIPTION copy_data` option.

Options

pg_createsubscriber accepts the following command-line arguments:

`-d dbname`

`--database=dbname`

The name of the database in which to create a subscription. Multiple databases can be selected by writing multiple `-d` switches.

`-D directory`

`--pgdata=directory`

The target directory that contains a cluster directory from a physical replica.

`-n`

`--dry-run`

Do everything except actually modifying the target directory.

`-p port`

`--subscriber-port=port`

The port number on which the target server is listening for connections. Defaults to running the target server on port 50432 to avoid unintended client connections.

`-P connstr`

`--publisher-server=connstr`

The connection string to the publisher. For details see [Section 35.1.1](#).

`-s dir`

`--socketdir=dir`

The directory to use for postmaster sockets on target server. The default is current directory.

`-t seconds`

`--recovery-timeout=seconds`

The maximum number of seconds to wait for recovery to end. Setting to 0 disables. The default is 0.

`-U username`

`--subscriber-username=username`

The user name to connect as on target server. Defaults to the current operating system user name.

`-v`

`--verbose`

Enables verbose mode. This will cause `pg_createsubscriber` to output progress messages and detailed information about each step to standard error. Repeating the option causes additional debug-level messages to appear on standard error.

`--config-file=filename`

Use the specified main server configuration file for the target data directory. `pg_createsubscriber` internally uses the `pg_ctl` command to start and stop the target server. It allows you to specify the actual `postgresql.conf` configuration file if it is stored outside the data directory.

`--publication=name`

The publication name to set up the logical replication. Multiple publications can be specified by writing multiple `--publication` switches. The number of publication names must match the number of specified databases, otherwise an error is reported. The order of the multiple publication name switches must match the order of database switches. If this option is not specified, a generated name is assigned to the publication name.

`--replication-slot=name`

The replication slot name to set up the logical replication. Multiple replication slots can be specified by writing multiple `--replication-slot` switches. The number of replication slot names must match the number of specified databases, otherwise an error is reported. The order of the multiple replication slot name switches must match the order of database switches. If this option is not specified, the subscription name is assigned to the replication slot name.

`--subscription=name`

The subscription name to set up the logical replication. Multiple subscriptions can be specified by writing multiple `--subscription` switches. The number of subscription names must match the number of specified databases, otherwise an error is reported. The order of the multiple subscription name switches must match the order of database switches. If this option is not specified, a generated name is assigned to the subscription name.

`-V`

`--version`

Print the `pg_createsubscriber` version and exit.

`-?`

`--help`

Show help about `pg_createsubscriber` command line arguments, and exit.

Notes

Prerequisites

There are some prerequisites for `pg_createsubscriber` to convert the target server into a logical replica. If these are not met, an error will be reported. The source and target servers must have the same major version as the `pg_createsubscriber`. The given target data directory must have the same system identifier

as the source data directory. The given database user for the target data directory must have privileges for creating `subscriptions` and using `pg_replication_origin_advance()`.

The target server must be used as a physical standby. The target server must have `max_replication_slots` and `max_logical_replication_workers` configured to a value greater than or equal to the number of specified databases. The target server must have `max_worker_processes` configured to a value greater than the number of specified databases. The target server must accept local connections.

The source server must accept connections from the target server. The source server must not be in recovery. The source server must have `wal_level` as `logical`. The source server must have `max_replication_slots` configured to a value greater than or equal to the number of specified databases plus existing replication slots. The source server must have `max_wal_senders` configured to a value greater than or equal to the number of specified databases and existing WAL sender processes.

Warnings

If `pg_createsubscriber` fails after the target server was promoted, then the data directory is likely not in a state that can be recovered. In such case, creating a new standby server is recommended.

`pg_createsubscriber` usually starts the target server with different connection settings during transformation. Hence, connections to the target server should fail.

Since DDL commands are not replicated by logical replication, avoid executing DDL commands that change the database schema while running `pg_createsubscriber`. If the target server has already been converted to logical replica, the DDL commands might not be replicated, which might cause an error.

If `pg_createsubscriber` fails while processing, objects (publications, replication slots) created on the source server are removed. The removal might fail if the target server cannot connect to the source server. In such a case, a warning message will inform the objects left. If the target server is running, it will be stopped.

If the replication is using `primary_slot_name`, it will be removed from the source server after the logical replication setup.

If the target server is a synchronous replica, transaction commits on the primary might wait for replication while running `pg_createsubscriber`.

`pg_createsubscriber` sets up logical replication with two-phase commit disabled. This means that any prepared transactions will be replicated at the time of `COMMIT PREPARED`, without advance preparation. Once setup is complete, you can manually drop and re-create the subscription(s) with the `two_phase` option enabled.

`pg_createsubscriber` changes the system identifier using `pg_resetwal`. It would avoid situations in which the target server might use WAL files from the source server. If the target server has a standby, replication will break and a fresh standby should be created.

How It Works

The basic idea is to have a replication start point from the source server and set up a logical replication to start from this point:

1. Start the target server with the specified command-line options. If the target server is already running, `pg_createsubscriber` will terminate with an error.
2. Check if the target server can be converted. There are also a few checks on the source server. If any of the prerequisites are not met, `pg_createsubscriber` will terminate with an error.
3. Create a publication and replication slot for each specified database on the source server. Each publication is created using `FOR ALL TABLES`. If the `--publication` option is not specified, the publication has the following name pattern: `"pg_createsubscriber_%u_%x"` (parameter: database *oid*, random *int*). If the `--replication-slot` option is not specified, the replication slot has the following name pattern: `"pg_createsubscriber_%u_%x"` (parameters: database *oid*, random *int*).

These replication slots will be used by the subscriptions in a future step. The last replication slot LSN is used as a stopping point in the [recovery_target_lsn](#) parameter and by the subscriptions as a replication start point. It guarantees that no transaction will be lost.

4. Write recovery parameters into the target data directory and restart the target server. It specifies an LSN ([recovery_target_lsn](#)) of the write-ahead log location up to which recovery will proceed. It also specifies `promote` as the action that the server should take once the recovery target is reached. Additional [recovery parameters](#) are added to avoid unexpected behavior during the recovery process such as end of the recovery as soon as a consistent state is reached (WAL should be applied until the replication start location) and multiple recovery targets that can cause a failure. This step finishes once the server ends standby mode and is accepting read-write transactions. If `--recovery-timeout` option is set, `pg_createsubscriber` terminates if recovery does not end until the given number of seconds.
5. Create a subscription for each specified database on the target server. If the `--subscription` option is not specified, the subscription has the following name pattern: `"pg_createsubscriber_%u_%x"` (parameters: database *oid*, random *int*). It does not copy existing data from the source server. It does not create a replication slot. Instead, it uses the replication slot that was created in a previous step. The subscription is created but it is not enabled yet. The reason is the replication progress must be set to the replication start point before starting the replication.
6. Drop publications on the target server that were replicated because they were created before the replication start location. It has no use on the subscriber.
7. Set the replication progress to the replication start point for each subscription. When the target server starts the recovery process, it catches up to the replication start point. This is the exact LSN to be used as a initial replication location for each subscription. The replication origin name is obtained since the subscription was created. The replication origin name and the replication start point are used in `pg_replication_origin_advance()` to set up the initial replication location.
8. Enable the subscription for each specified database on the target server. The subscription starts applying transactions from the replication start point.
9. If the standby server was using [primary_slot_name](#), it has no use from now on so drop it.
10. If the standby server contains [failover replication slots](#), they cannot be synchronized anymore, so drop them.
11. Update the system identifier on the target server. The [pg_resetwal](#) is run to modify the system identifier. The target server is stopped as a `pg_resetwal` requirement.

Examples

To create a logical replica for databases `hr` and `finance` from a physical replica at `foo`:

```
$ pg_createsubscriber -D /usr/local/pgsql/data -P "host=foo" -d hr -d finance
```

pg_ctl

pg_ctl — initialize, start, stop, or control a Postgres Pro server

Synopsis

```
pg_ctl init[db] [-D datadir] [-s] [-o initdb-options]  
  
pg_ctl start [-D datadir] [-l filename] [-W] [-t seconds] [-s] [-o options] [-p path] [-c]  
  
pg_ctl stop [-D datadir] [-m s[mart] | f[ast] | i[mmediate]] [-W] [-t seconds] [-s]  
  
pg_ctl restart [-D datadir] [-m s[mart] | f[ast] | i[mmediate]] [-W] [-t seconds] [-s] [-o options]  
[-c]  
  
pg_ctl reload [-D datadir] [-s]  
  
pg_ctl status [-D datadir]  
  
pg_ctl promote [-D datadir] [-W] [-t seconds] [-s]  
  
pg_ctl logrotate [-D datadir] [-s]  
  
pg_ctl kill signal_name process_id
```

On Microsoft Windows, also:

```
pg_ctl register [-D datadir] [-N servicename] [-U username] [-P password] [-S a[uto] | d[emand]] [-e source] [-W] [-t seconds] [-s] [-o options]  
  
pg_ctl unregister [-N servicename]
```

Description

pg_ctl is a utility for initializing a Postgres Pro database cluster, starting, stopping, or restarting the Postgres Pro database server ([postgres](#)), or displaying the status of a running server. Although the server can be started manually, pg_ctl encapsulates tasks such as redirecting log output and properly detaching from the terminal and process group. It also provides convenient options for controlled shutdown.

The `init` or `initdb` mode creates a new Postgres Pro database cluster, that is, a collection of databases that will be managed by a single server instance. This mode invokes the `initdb` command. See [initdb](#) for details.

`start` mode launches a new server. The server is started in the background, and its standard input is attached to `/dev/null` (or `nul` on Windows). On Unix-like systems, by default, the server's standard output and standard error are sent to pg_ctl's standard output (not standard error). The standard output of pg_ctl should then be redirected to a file or piped to another process such as a log rotating program like `rotatelogs`; otherwise `postgres` will write its output to the controlling terminal (from the background) and will not leave the shell's process group. On Windows, by default the server's standard output and standard error are sent to the terminal. These default behaviors can be changed by using `-l` to append the server's output to a log file. Use of either `-l` or output redirection is recommended.

`stop` mode shuts down the server that is running in the specified data directory. Three different shutdown methods can be selected with the `-m` option. “Smart” mode disallows new connections, then waits for all existing clients to disconnect. If the server is in hot standby, recovery and streaming replication will be terminated once all clients have disconnected. “Fast” mode (the default) does not wait for clients to disconnect. All active transactions are rolled back and clients are forcibly disconnected, then the server is shut down. “Immediate” mode will abort all server processes immediately, without a clean shutdown. This choice will lead to a crash-recovery cycle during the next server start.

restart mode effectively executes a stop followed by a start. This allows changing the `postgres` command-line options, or changing configuration-file options that cannot be changed without restarting the server. If relative paths were used on the command line during server start, `restart` might fail unless `pg_ctl` is executed in the same current directory as it was during server start.

reload mode simply sends the `postgres` server process a `SIGHUP` signal, causing it to reread its configuration files (`postgresql.conf`, `pg_hba.conf`, etc.). This allows changing configuration-file options that do not require a full server restart to take effect.

status mode checks whether a server is running in the specified data directory. If it is, the server's PID and the command line options that were used to invoke it are displayed. If the server is not running, `pg_ctl` returns an exit status of 3. If an accessible data directory is not specified, `pg_ctl` returns an exit status of 4.

promote mode commands the standby server that is running in the specified data directory to end standby mode and begin read-write operations.

logrotate mode rotates the server log file. For details on how to use this mode with external log rotation tools, see [Section 24.4](#).

kill mode sends a signal to a specified process. This is primarily valuable on Microsoft Windows which does not have a built-in kill command. Use `--help` to see a list of supported signal names.

register mode registers the Postgres Pro server as a system service on Microsoft Windows. The `-s` option allows selection of service start type, either “auto” (start service automatically on system startup) or “demand” (start service on demand).

unregister mode unregisters a system service on Microsoft Windows. This undoes the effects of the `register` command.

Options

`-c`
`--core-files`

Attempt to allow server crashes to produce core files, on platforms where this is possible, by lifting any soft resource limit placed on core files. This is useful in debugging or diagnosing problems by allowing a stack trace to be obtained from a failed server process.

`-D datadir`
`--pgdata=datadir`

Specifies the file system location of the database configuration files. If this option is omitted, the environment variable `PGDATA` is used.

`-l filename`
`--log=filename`

Append the server log output to *filename*. If the file does not exist, it is created. The umask is set to 077, so access to the log file is disallowed to other users by default.

`-m mode`
`--mode=mode`

Specifies the shutdown mode. *mode* can be `smart`, `fast`, or `immediate`, or the first letter of one of these three. If this option is omitted, `fast` is the default.

`-o options`
`--options=options`

Specifies options to be passed directly to the `postgres` command. `-o` can be specified multiple times, with all the given options being passed through.

The *options* should usually be surrounded by single or double quotes to ensure that they are passed through as a group.

`-o initdb-options`
`--options=initdb-options`

Specifies options to be passed directly to the `initdb` command. `-o` can be specified multiple times, with all the given options being passed through.

The *initdb-options* should usually be surrounded by single or double quotes to ensure that they are passed through as a group.

`-p path`

Specifies the location of the `postgres` executable. By default the `postgres` executable is taken from the same directory as `pg_ctl`, or failing that, the hard-wired installation directory. It is not necessary to use this option unless you are doing something unusual and get errors that the `postgres` executable was not found.

In `init` mode, this option analogously specifies the location of the `initdb` executable.

`-s`
`--silent`

Print only errors, no informational messages.

`-t seconds`
`--timeout=seconds`

Specifies the maximum number of seconds to wait when waiting for an operation to complete (see option `-w`). Defaults to the value of the `PGCTLTIMEOUT` environment variable or, if not set, to 60 seconds.

`-V`
`--version`

Print the `pg_ctl` version and exit.

`-w`
`--wait`

Wait for the operation to complete. This is supported for the modes `start`, `stop`, `restart`, `promote`, and `register`, and is the default for those modes.

When waiting, `pg_ctl` repeatedly checks the server's PID file, sleeping for a short amount of time between checks. Startup is considered complete when the PID file indicates that the server is ready to accept connections. Shutdown is considered complete when the server removes the PID file. `pg_ctl` returns an exit code based on the success of the startup or shutdown.

If the operation does not complete within the timeout (see option `-t`), then `pg_ctl` exits with a nonzero exit status. But note that the operation might continue in the background and eventually succeed.

`-W`
`--no-wait`

Do not wait for the operation to complete. This is the opposite of the option `-w`.

If waiting is disabled, the requested action is triggered, but there is no feedback about its success. In that case, the server log file or an external monitoring system would have to be used to check the progress and success of the operation.

In prior releases of Postgres Pro, this was the default except for the `stop` mode.

-?
--help

Show help about pg_ctl command line arguments, and exit.

If an option is specified that is valid, but not relevant to the selected operating mode, pg_ctl ignores it.

Options for Windows

-e *source*

Name of the event source for pg_ctl to use for logging to the event log when running as a Windows service. The default is `Postgres Pro`. Note that this only controls messages sent from pg_ctl itself; once started, the server will use the event source specified by its [event_source](#) parameter. Should the server fail very early in startup, before that parameter has been set, it might also log using the default event source name `Postgres Pro`.

-N *servicename*

Name of the system service to register. This name will be used as both the service name and the display name. The default is `Postgres Pro`.

-P *password*

Password for the user to run the service as.

-S *start-type*

Start type of the system service. *start-type* can be `auto`, or `demand`, or the first letter of one of these two. If this option is omitted, `auto` is the default.

-U *username*

User name for the user to run the service as. For domain users, use the format `DOMAIN\username`.

Environment

PGCTLTIMEOUT

Default limit on the number of seconds to wait when waiting for startup or shutdown to complete. If not set, the default is 60 seconds.

PGDATA

Default data directory location.

Most pg_ctl modes require knowing the data directory location; therefore, the `-D` option is required unless PGDATA is set.

For additional variables that affect the server, see [postgres](#).

Files

postmaster.pid

pg_ctl examines this file in the data directory to determine whether the server is currently running.

postmaster.opts

If this file exists in the data directory, pg_ctl (in `restart` mode) will pass the contents of the file as options to postgres, unless overridden by the `-o` option. The contents of this file are also displayed in `status` mode.

Examples

Starting the Server

To start the server, waiting until the server is accepting connections:


```
$ pg_ctl start
```

To start the server using port 5433, and running without `fsync`, use:

```
$ pg_ctl -o "-F -p 5433" start
```

Stopping the Server

To stop the server, use:

```
$ pg_ctl stop
```

The `-m` option allows control over *how* the server shuts down:

```
$ pg_ctl stop -m smart
```

Restarting the Server

Restarting the server is almost equivalent to stopping the server and starting it again, except that by default, `pg_ctl` saves and reuses the command line options that were passed to the previously-running instance. To restart the server using the same options as before, use:

```
$ pg_ctl restart
```

But if `-o` is specified, that replaces any previous options. To restart using port 5433, disabling `fsync` upon restart:

```
$ pg_ctl -o "-F -p 5433" restart
```

Showing the Server Status

Here is sample status output from `pg_ctl`:

```
$ pg_ctl status
```

```
pg_ctl: server is running (PID: 13718)
/usr/local/pgsql/bin/postgres "-D" "/usr/local/pgsql/data" "-p" "5433" "-B" "128"
```

The second line is the command that would be invoked in restart mode.

See Also

[initdb](#), [postgres](#)

pgpro_tune

pgpro_tune — a command-line tool for automatic tuning

Synopsis

```
pgpro_tune [option...] [ -D | --pgdata ] directory preset_name
```

Description

The pgpro_tune utility is a command-line tool for automatic tuning. Optimal values for different configuration parameters of Postgres Pro Shardman depend on the hardware configuration. The pgpro_tune utility collects information about the system and transforms it into a set of parameters written into the configuration file. Note that it must be aware of the Shardman server topology to work properly.

Usage

The use of pgpro_tune is enabled in initdb by default. There are two initdb options for this tool: `--no-tune` disables its usage, `--tune=OPTIONS` allows passing additional options for the pgpro_tune run.

First pgpro_tune runs the default preset. Additional presets may be provided by using `-P` or `--preset=NAME` option, or specifying the preset name as the last command-line argument. To disable the default preset, use the `--no-default` option.

Additional options may be passed via `-O` or `--options=OPTIONS` parameter for any preset, but for the default preset `--default-options=OPTIONS` should be used.

Additional environment variables can be set for presets with the `--set NAME=VALUE` parameter. These values will be the same both for the default preset and all custom presets.

By default, the presets are placed in the `share` directory, but you can specify their location manually in the `--preset-dir=NAME` parameter.

By default, pgpro_tune works with the `postgresql.conf` configuration file located in the data directory. The data directory can be specified with `-D` or `--pgdata=DATADIR`, or provided with environment variable `PGDATA`. When running pgpro_tune from initdb, the tool will use the data directory of initdb.

Configuration file can also be specified with `--config-file=FILENAME`.

The pgpro_tune tool writes its changes into a separate block at the end of the configuration file, which is commented as added by pgpro_tune. For details, see [the example of the preset](#).

Options

`--config-file=filename`

Specifies the main configuration file name.

`-D datadir`

`--pgdata=datadir`

Specifies the directory where the database cluster should be stored.

`--default-options=options`

Specifies parameters for the default preset.

`--no-default`

Disables the default preset. Note that in this case another preset must be specified anyway either with the `-P` option or as the last parameter with the preset name.

`-O preset-options`

`--options=preset-options`

Specifies options to be passed directly to the preset that was previously specified in the `-P` parameter or to the preset indicated as the last command-line parameter, if `-P` was not used.

`-P name`

`--preset=name`

Specifies the preset name. If several names are indicated, presets are executed in turn, taking into account the output of the previous preset.

Note

The preset name may be specified either with this option or as the last parameter. Using both will cause an error.

`--preset-dir=name`

Specifies the directory containing presets. By default, presets are located in the `share` directory.

`--set name=value`

Sets the environment variable for presets.

`--show`

Shows all available presets from the directory specified in the `--preset-dir` option, if it is set, or from the default directory otherwise.

`-V`

`--version`

Print the `pgpro_tune` version and exit.

`-?`

`--help`

Show help about `pgpro_tune` command line arguments, and exit.

Environment variables

`pgpro_tune` sets the following environment variables before running presets.

- `EDITION` — the Postgres Pro edition.
- `ENABLE_CRASH_INFO` — if set, Postgres Pro was built with the `--enable-crash-info` option.
- `ENABLE_NLS` — if set, Postgres Pro was built with the `--enable-nls` option.
- `ENABLE_PGPRO_TUNE` — if set, Postgres Pro was built with the `--enable-pgpro-tune` option.
- `MEMMB` — the RAM size in megabytes.
- `MVER` — Postgres Pro major version.
- `NCPU` — the number of CPUs.
- `USE_BONJOUR` — if set, Postgres Pro was built with the `--with-bonjour` option.
- `USE_BSD_AUTH` — if set, Postgres Pro was built with the `--with-bsd-auth` option.
- `USE_ICU` — if set, Postgres Pro was built the `--with-icu` option.
- `USE_LDAP` — if set, Postgres Pro was built with the `--with-ldap` option.
- `USE_LIBUNWIND` — if set, Postgres Pro was built with the `--with-libunwind` option.

- `USE_LIBXML` — if set, Postgres Pro was built with the `--with-libxml` option.
- `USE_LIBXSLT` — if set, Postgres Pro was built with the `--with-libxslt` option.
- `USE_LLVM` — if set, Postgres Pro was built with the `--with-llvm` option.
- `USE_LZ4` — if set, Postgres Pro was built with the `--with-lz4` option.
- `USE_OPENSSL` — if set, Postgres Pro was built with the `--with-openssl` option.
- `USE_PAM` — if set, Postgres Pro was built with the `--with-pam` option.
- `USE_SYSTEMD` — if set, Postgres Pro was built with the `--with-systemd` option.
- `USE_ZSTD` — if set, Postgres Pro was built with the `--with-zstd` option.

Writing presets

A preset is an executable, usually a shell script, that transforms the information passed through environment variables into the set of configuration parameters.

The expected output of the preset is the set of specifically formatted lines. The first non-whitespace character is the indicator of the line content interpretation for `pgpro_tune`. The indicator should be one of the following:

- `#` for the comment action, the rest of the output line will be written into the configuration file after `'#'`.
- `=` for the replace action, the rest of the output line should contain a pair of parameters `NAME` and `VALUE` (both `NAME=VALUE` and `NAME VALUE` are acceptable). Line `NAME = VALUE` will be written into the configuration file. The text after `VALUE` will be added as a comment.
- `+=` or `+=` for the add action, the rest of the output line should contain a pair of parameters `NAME` and `VALUE`. `VALUE` will be added to the previous effective value of the parameter `NAME` at the back (`+=`) or at the front (`+=`). If there is no effective value, the `ADD` action works like the replace action. The text after `VALUE` will be added as a comment.
- `--` for the withdraw action, the rest of the output line should contain a pair of parameters `NAME` and `VALUE`. If the `NAME` parameter has the effective value containing `VALUE`, it will be removed and the rest of the effective value will remain the same. The text after `VALUE` will be added as a comment. If there is no effective value, nothing happens.

Example of the preset

Suppose you have the following configuration file called `test.conf`:

```
work_mem = 4MB # Default value
shared_preload_libraries = 'plantuner'
search_path = '"$user",wrong_schema,public'
```

And the following preset called `test.tune`:

```
echo "# Adding new configuration parameters."
#Replace configuration parameter value by a new one
echo "work_mem = 8MB"
#Append to the start of existing value
echo "shared_preload_libraries += pg_stat_statements"
#Append to the end of existing value
echo "shared_preload_libraries += pg_prewarm"
#Withdraw from existing value
echo "search_path -= 'wrong_schema'"
```

Run the following command to use this preset:

```
pgpro_tune --config-file=/path/to/test.conf -P/path/to/test.tune --no-default
```

This command will result in the following changes of the configuration file:

```
#-----
# The following settings were added by pgpro_tune.
# pgpro_tune was run with the following options:
# --no-default --config-file=/path/to/test/conf -P/path/to/test/tune --no-default
# At YYYY-MM-DD HH:MM:SS
#-----
# Adding new configuration parameters.
work_mem = 8MB
shared_preload_libraries = 'pg_stat_statements, plantuner'
shared_preload_libraries = 'pg_stat_statements, plantuner, pg_prewarm'
search_path = '"$user", public'
#-----
# End of settings added by pgpro_tune at YYYY-MM-DD HH-MM-SS
#-----
```

pg_resetwal

`pg_resetwal` — reset the write-ahead log and other control information of a Postgres Pro database cluster

Synopsis

```
pg_resetwal [ -f | --force ] [ -n | --dry-run ] [option...] [ -D | --pgdata ] datadir
```

Description

`pg_resetwal` clears the write-ahead log (WAL) and optionally resets some other control information stored in the `pg_control` file. This function is sometimes needed if these files have become corrupted. It should be used only as a last resort, when the server will not start due to such corruption.

Some options, such as `--wal-segsize` (see below), can also be used to modify certain global settings of a database cluster without the need to rerun `initdb`. This can be done safely on an otherwise sound database cluster, if none of the dangerous modes mentioned below are used.

If `pg_resetwal` is used on a data directory where the server has been cleanly shut down and the control file is sound, then it will have no effect on the contents of the database system, except that no longer used WAL files are cleared away. Any other use is potentially dangerous and must be done with great care. `pg_resetwal` will require the `-f` (force) option to be specified before working on a data directory in an unclean shutdown state or with a corrupted control file.

After running this command on a data directory with corrupted WAL or a corrupted control file, it should be possible to start the server, but bear in mind that the database might contain inconsistent data due to partially-committed transactions. You should immediately dump your data, run `initdb`, and restore. After restore, check for inconsistencies and repair as needed.

If `pg_resetwal` complains that it cannot determine valid data for `pg_control`, you can force it to proceed anyway by specifying the `-f` (force) option. In this case plausible values will be substituted for the missing data. Most of the fields can be expected to match, but manual assistance might be needed for the next OID, next transaction ID, next multitransaction ID and offset, and WAL starting location fields. These fields can be set using the options discussed below. If you are not able to determine correct values for all these fields, `-f` can still be used, but the recovered database must be treated with even more suspicion than usual: an immediate dump and restore is imperative. *Do not* execute any data-modifying operations in the database before you dump, as any such action is likely to make the corruption worse.

This utility can only be run by the user who installed the server, because it requires read/write access to the data directory.

Options

```
datadir
-D datadir
--pgdata=datadir
```

Specifies the location of the database directory. For safety reasons, you must specify the data directory on the command line. `pg_resetwal` does not use the environment variable `PGDATA`.

```
-f
--force
```

Force `pg_resetwal` to proceed even in situations where it could be dangerous, as explained above. Specifically, this option is required to proceed if the server had not been cleanly shut down or if `pg_resetwal` cannot determine valid data for `pg_control`.

`-n`
`--dry-run`

The `-n/--dry-run` option instructs `pg_resetwal` to print the values reconstructed from `pg_control` and values about to be changed, and then exit without modifying anything. This is mainly a debugging tool, but can be useful as a sanity check before allowing `pg_resetwal` to proceed for real.

`-V`
`--version`

Display version information, then exit.

`-?`
`--help`

Show help, then exit.

The following options are only needed when `pg_resetwal` is unable to determine appropriate values by reading `pg_control`. Safe values can be determined as described below. For values that take numeric arguments, hexadecimal values can be specified by using the prefix `0x`. Note that these instructions only apply with the standard block size of 8 kB.

`-c xid,xid`
`--commit-timestamp-ids=xid,xid`

Manually set the oldest and newest transaction IDs for which the commit time can be retrieved.

A safe value for the oldest transaction ID for which the commit time can be retrieved (first part) can be determined by looking for the numerically smallest file name in the directory `pg_commit_ts` under the data directory. Conversely, a safe value for the newest transaction ID for which the commit time can be retrieved (second part) can be determined by looking for the numerically greatest file name in the same directory. The file names are in hexadecimal.

`-l walfile`
`--next-wal-file=walfile`

Manually set the WAL starting location by specifying the name of the next WAL segment file.

The name of next WAL segment file should be larger than any WAL segment file name currently existing in the directory `pg_wal` under the data directory. These names are also in hexadecimal and have three parts. The first part is the “timeline ID” and should usually be kept the same. For example, if `000000010000000320000004A` is the largest entry in `pg_wal`, use `-l 000000010000000320000004B` or higher.

Note that when using nondefault WAL segment sizes, the numbers in the WAL file names are different from the LSNs that are reported by system functions and system views. This option takes a WAL file name, not an LSN.

Note

`pg_resetwal` itself looks at the files in `pg_wal` and chooses a default `-l` setting beyond the last existing file name. Therefore, manual adjustment of `-l` should only be needed if you are aware of WAL segment files that are not currently present in `pg_wal`, such as entries in an offline archive; or if the contents of `pg_wal` have been lost entirely.

`-m mxid,mxid`
`--multixact-ids=mxid,mxid`

Manually set the next and oldest multitransaction ID.

A safe value for the next multitransaction ID (first part) can be determined by looking for the numerically largest file name in the directory `pg_multixact/offsets` under the data directory, adding one, and then multiplying by 2097152 (0x200000). Conversely, a safe value for the oldest multitransaction ID (second part of `-m`) can be determined by looking for the numerically smallest file name in the same directory and multiplying by 2097152. The file names are in hexadecimal, so the easiest way to do this is to specify the option value in hexadecimal and append four zeroes.

```
-o oid
--next-oid=oid
```

Manually set the next OID.

There is no comparably easy way to determine a next OID that's beyond the largest one in the database, but fortunately it is not critical to get the next-OID setting right.

```
-O mxoff
--multixact-offset=mxoff
```

Manually set the next multitransaction offset.

A safe value can be determined by looking for the numerically largest file name in the directory `pg_multixact/members` under the data directory, adding one, and then multiplying by 1851392 (0x1c4000). The file names are in hexadecimal. There is no simple recipe such as the ones for other options of appending zeroes.

```
--wal-segsize=wal_segment_size
```

Set the new WAL segment size, in megabytes. The value must be set to a power of 2 between 1 and 1024 (megabytes). See the same option of [initdb](#) for more information.

This option can also be used to change the WAL segment size of an existing database cluster, avoiding the need to `re-initdb`.

Note

While `pg_resetwal` will set the WAL starting address beyond the latest existing WAL segment file, some segment size changes can cause previous WAL file names to be reused. It is recommended to use `-l` together with this option to manually set the WAL starting address if WAL file name overlap will cause problems with your archiving strategy.

```
-u xid
--oldest-transaction-id=xid
```

Manually set the oldest unfrozen transaction ID.

A safe value can be determined by looking for the numerically smallest file name in the directory `pg_xact` under the data directory and then multiplying by 67108864 (0x4000000). Note that the file names are in hexadecimal. It is usually easiest to specify the option value in hexadecimal too. For example, if `0007` is the smallest entry in `pg_xact`, `-u 0x700000` will work (five trailing zeroes provide the proper multiplier).

```
-x xid
--next-transaction-id=xid
```

Manually set the next transaction ID.

A safe value can be determined by looking for the numerically largest file name in the directory `pg_xact` under the data directory, adding one, and then multiplying by 67108864 (0x4000000). Note that the file names are in hexadecimal. It is usually easiest to specify the option value in hexadecimal too. For example, if `0011` is the largest entry in `pg_xact`, `-x 0x1200000` will work (five trailing zeroes provide the proper multiplier).

Environment

PG_COLOR

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

This command must not be used when the server is running. `pg_resetwal` will refuse to start up if it finds a server lock file in the data directory. If the server crashed then a lock file might have been left behind; in that case you can remove the lock file to allow `pg_resetwal` to run. But before you do so, make doubly certain that there is no server process still alive.

`pg_resetwal` works only with servers of the same major version.

pg_rewind

`pg_rewind` — synchronize a Postgres Pro data directory with another data directory that was forked from it

Synopsis

```
pg_rewind [option...] { -D | --target-pgdata } directory { --source-pgdata=directory | --source-server=connstr }
```

Description

`pg_rewind` is a tool for synchronizing a Postgres Pro cluster with another copy of the same cluster, after the clusters' timelines have diverged. A typical scenario is to bring an old primary server back online after failover as a standby that follows the new primary.

After a successful rewind, the state of the target data directory is analogous to a base backup of the source data directory. Unlike taking a new base backup or using a tool like `rsync`, `pg_rewind` does not require comparing or copying unchanged relation blocks in the cluster. Only changed blocks from existing relation files are copied; all other files, including new relation files, configuration files, and WAL segments, are copied in full. As such the rewind operation is significantly faster than other approaches when the database is large and only a small fraction of blocks differ between the clusters.

`pg_rewind` examines the timeline histories of the source and target clusters to determine the point where they diverged, and expects to find WAL in the target cluster's `pg_wal` directory reaching all the way back to the point of divergence. The point of divergence can be found either on the target timeline, the source timeline, or their common ancestor. In the typical failover scenario where the target cluster was shut down soon after the divergence, this is not a problem, but if the target cluster ran for a long time after the divergence, its old WAL files might no longer be present. In this case, you can manually copy them from the WAL archive to the `pg_wal` directory, or run `pg_rewind` with the `-c` option to automatically retrieve them from the WAL archive. The use of `pg_rewind` is not limited to failover, e.g., a standby server can be promoted, run some write transactions, and then rewound to become a standby again.

After running `pg_rewind`, WAL replay needs to complete for the data directory to be in a consistent state. When the target server is started again it will enter archive recovery and replay all WAL generated in the source server from the last checkpoint before the point of divergence. If some of the WAL was no longer available in the source server when `pg_rewind` was run, and therefore could not be copied by the `pg_rewind` session, it must be made available when the target server is started. This can be done by creating a `recovery.signal` file in the target data directory and by configuring a suitable [restore_command](#) in `postgresql.conf`.

`pg_rewind` requires that the target server either has the [wal_log_hints](#) option enabled in `postgresql.conf` or data checksums enabled when the cluster was initialized with `initdb`. While `wal_log_hints` is currently off by default, the checksums are enabled. [full_page_writes](#) must also be set to `on`, and it is enabled by default.

Warning: Failures While Rewinding

If `pg_rewind` fails while processing, then the data folder of the target is likely not in a state that can be recovered. In such a case, taking a new fresh backup is recommended.

As `pg_rewind` copies configuration files entirely from the source, it may be required to correct the configuration used for recovery before restarting the target server, especially if the target is reintroduced as a standby of the source. If you restart the server after the rewind operation has finished but without configuring recovery, the target may again diverge from the primary.

`pg_rewind` will fail immediately if it finds files it cannot write directly to. This can happen for example when the source and the target server use the same file mapping for read-only SSL keys

and certificates. If such files are present on the target server it is recommended to remove them before running `pg_rewind`. After doing the rewind, some of those files may have been copied from the source, in which case it may be necessary to remove the data copied and restore back the set of links used before the rewind.

Options

`pg_rewind` accepts the following command-line arguments:

`-D directory`

`--target-pgdata=directory`

This option specifies the target data directory that is synchronized with the source. The target server must be shut down cleanly before running `pg_rewind`

`--source-pgdata=directory`

Specifies the file system path to the data directory of the source server to synchronize the target with. This option requires the source server to be cleanly shut down.

`--source-server=connstr`

Specifies a libpq connection string to connect to the source Postgres Pro server to synchronize the target with. The connection must be a normal (non-replication) connection with a role having sufficient permissions to execute the functions used by `pg_rewind` on the source server (see Notes section for details) or a superuser role. This option requires the source server to be running and accepting connections.

`-R`

`--write-recovery-conf`

Create `standby.signal` and append connection settings to `postgresql.auto.conf` in the output directory. `--source-server` is mandatory with this option.

`-n`

`--dry-run`

Do everything except actually modifying the target directory.

`-N`

`--no-sync`

By default, `pg_rewind` will wait for all files to be written safely to disk. This option causes `pg_rewind` to return without waiting, which is faster, but means that a subsequent operating system crash can leave the data directory corrupt. Generally, this option is useful for testing but should not be used on a production installation.

`-P`

`--progress`

Enables progress reporting. Turning this on will deliver an approximate progress report while copying data from the source cluster.

`-c`

`--restore-target-wal`

Use `restore_command` defined in the target cluster configuration to retrieve WAL files from the WAL archive if these files are no longer available in the `pg_wal` directory.

`--config-file=filename`

Use the specified main server configuration file for the target cluster. This affects `pg_rewind` when it uses internally the `postgres` command for the rewind operation on this cluster (when retrieving

`restore_command` with the option `-c/--restore-target-wal` and when forcing a completion of crash recovery).

`--debug`

Print verbose debugging output that is mostly useful for developers debugging `pg_rewind`.

`--no-ensure-shutdown`

`pg_rewind` requires that the target server is cleanly shut down before rewinding. By default, if the target server is not shut down cleanly, `pg_rewind` starts the target server in single-user mode to complete crash recovery first, and stops it. By passing this option, `pg_rewind` skips this and errors out immediately if the server is not cleanly shut down. Users are expected to handle the situation themselves in that case.

`--sync-method=method`

When set to `fsync`, which is the default, `pg_rewind` will recursively open and synchronize all files in the data directory. The search for files will follow symbolic links for the WAL directory and each configured tablespace.

On Linux, `syncfs` may be used instead to ask the operating system to synchronize the whole file systems that contain the data directory, the WAL files, and each tablespace. See [recovery_init_sync_method](#) for information about the caveats to be aware of when using `syncfs`.

This option has no effect when `--no-sync` is used.

`--biha`

Set the `biha` node to the `CSTATE_FORMING` state for the subsequent status of the node to be set automatically depending on the current status of the cluster. This option is used when restoring the node of the built-in high-availability cluster from the `NODE_ERROR` state.

Important

The `--biha` option is essential for saving `biha` configuration files. Using `pg_rewind` in a BiHA cluster without the `--biha` option may cause cluster configuration inconsistency.

`-V`

`--version`

Display version information, then exit.

`-?`

`--help`

Show help, then exit.

Environment

When `--source-server` option is used, `pg_rewind` also uses the environment variables supported by `libpq` (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

When executing `pg_rewind` using an online cluster as source, a role having sufficient permissions to execute the functions used by `pg_rewind` on the source cluster can be used instead of a superuser. Here is how to create such a role, named `rewind_user` here:

```
CREATE USER rewind_user LOGIN;
GRANT EXECUTE ON function pg_catalog.pg_ls_dir(text, boolean, boolean) TO rewind_user;
GRANT EXECUTE ON function pg_catalog.pg_stat_file(text, boolean) TO rewind_user;
GRANT EXECUTE ON function pg_catalog.pg_read_binary_file(text) TO rewind_user;
GRANT EXECUTE ON function pg_catalog.pg_read_binary_file(text, bigint, bigint, boolean)
    TO rewind_user;
```

How It Works

The basic idea is to copy all file system-level changes from the source cluster to the target cluster:

1. Scan the WAL log of the target cluster, starting from the last checkpoint before the point where the source cluster's timeline history forked off from the target cluster. For each WAL record, record each data block that was touched. This yields a list of all the data blocks that were changed in the target cluster, after the source cluster forked off. If some of the WAL files are no longer available, try re-running `pg_rewind` with the `-c` option to search for the missing files in the WAL archive.
2. Copy all those changed blocks from the source cluster to the target cluster, either using direct file system access (`--source-pgdata`) or SQL (`--source-server`). Relation files are now in a state equivalent to the moment of the last completed checkpoint prior to the point at which the WAL timelines of the source and target diverged plus the current state on the source of any blocks changed on the target after that divergence.
3. Copy all other files, including new relation files, WAL segments, `pg_xact`, and configuration files from the source cluster to the target cluster. Similarly to base backups, the contents of the directories `pg_dynshmem/`, `pg_notify/`, `pg_replslot/`, `pg_serial/`, `pg_snapshots/`, `pg_stat_tmp/`, and `pg_subtrans/` are omitted from the data copied from the source cluster. The files `backup_label`, `tablespace_map`, `pg_internal.init`, `postmaster.opts`, `postmaster.pid` and `.DS_Store` as well as any file or directory beginning with `pgsql_tmp`, are omitted.
4. Create a `backup_label` file to begin WAL replay at the checkpoint created at failover and configure the `pg_control` file with a minimum consistency LSN defined as the result of `pg_current_wal_insert_lsn()` when rewinding from a live source or the last checkpoint LSN when rewinding from a stopped source.
5. When starting the target, Postgres Pro replays all the required WAL, resulting in a data directory in a consistent state.

pg-setup

pg-setup — set up a new Postgres Pro database cluster and manage the corresponding service

Synopsis

```
pg-setup initdb [initdb_options]
```

```
pg-setup find-free-port
```

```
pg-setup set-server-port port
```

```
pg-setup set parameter value
```

```
pg-setup service service_option
```

Description

pg-setup is a shell script provided in the Postgres Pro distribution to automate database cluster setup on Linux systems. This script is provided as part of the `postgrespro-sdm-17-server` package. Once Postgres Pro is installed, you can find pg-setup in the `install-dir/bin` directory, where `install-dir` is `/opt/pgpro/sdm-17`.

pg-setup must be run as root, but performs database administration operations as user `postgres`. You can run this script with different options to:

- initialize the database cluster
- configure the database cluster for a specific Postgres Pro product
- check for available ports and change the port used by Postgres Pro server
- enable/disable automatic startup of Postgres Pro service
- start, stop, or restart Postgres Pro service

Options

pg-setup accepts the following command-line arguments:

```
initdb [initdb_options]
```

Initialize the database cluster on behalf of the `postgres` user.

By default, the database cluster, configured for your Postgres Pro distribution, is initialized in the `/var/lib/pgpro/sdm-17/data` directory, with checksums enabled, `auth-local` parameter set to `peer`, and `auth-host` parameter set to `md5`. Localization settings are inherited from the `LANG` environment variable for the current session. All the `LC_*` environment variables are ignored. Optionally, you can provide `initdb` options to customize the installation.

If the default database is created using pg-setup, the path to its data directory is stored in the `/etc/default/postgrespro-shardman-17` file, so all the subsequent pg-setup commands, as well as any commands that manage Postgres Pro service, affect this database only. You cannot manage several databases using pg-setup.

For systems where more than one database server and/or application will run, you may need to adjust the configuration since pg-setup chooses the configuration settings depending on hardware characteristics, assuming that the system will use only one database server.

```
find-free-port
```

Search for a free port on your system. This option is useful if you are going to install more than one server instance, or the default 5432 port is used by another program.

`set-server-port port`

Specify the port number on which the server will listen for connections. Use this option to avoid conflicts if you are installing more than one server instance on the same system.

Default: 5432

`set name value`

Set the specified configuration parameter to the provided *value* in the `postgresql.conf` file. If this parameter has been already defined by the `ALTER SYSTEM` command, its previous value is removed from the `postgresql.auto.conf` file.

`service service_option`

Manage Postgres Pro service using one of the following options:

- `enable` — enable automatic service startup upon system restart.
- `disable` — disable automatic service startup upon system restart.
- `start` — start the service.
- `stop` — stop the service.
- `condrestart` — restart the service if it is running when `pg-setup` is invoked.
- `status` — return the Postgres Pro service status.

Notes

If you are installing Postgres Pro from the `postgrespro-sdm-17` package, `pg-setup` is invoked automatically with the default settings. As a result, the database cluster is initialized and the default database is created in the `/var/lib/pgpro/sdm-17/data` directory, Postgres Pro service autostart is enabled, and the service is started.

If you are installing Postgres Pro server directly from the `postgrespro-sdm-17-server` package, you can run this script manually to initialize the database cluster or manage the Postgres Pro service.

For details on binary installation specifics on Linux, see [Section 17.1](#).

pg_test_fsync

`pg_test_fsync` — determine fastest `wal_sync_method` for Postgres Pro

Synopsis

```
pg_test_fsync [option...]
```

Description

`pg_test_fsync` is intended to give you a reasonable idea of what the fastest [wal_sync_method](#) is on your specific system, as well as supplying diagnostic information in the event of an identified I/O problem. However, differences shown by `pg_test_fsync` might not make any significant difference in real database throughput, especially since many database servers are not speed-limited by their write-ahead logs. `pg_test_fsync` reports average file sync operation time in microseconds for each `wal_sync_method`, which can also be used to inform efforts to optimize the value of [commit_delay](#).

Options

`pg_test_fsync` accepts the following command-line options:

- `-f`
`--filename`
Specifies the file name to write test data in. This file should be in the same file system that the `pg_wal` directory is or will be placed in. (`pg_wal` contains the WAL files.) The default is `pg_test_fsync.out` in the current directory.
- `-s`
`--secs-per-test`
Specifies the number of seconds for each test. The more time per test, the greater the test's accuracy, but the longer it takes to run. The default is 5 seconds, which allows the program to complete in under 2 minutes.
- `-V`
`--version`
Print the `pg_test_fsync` version and exit.
- `-?`
`--help`
Show help about `pg_test_fsync` command line arguments, and exit.

Environment

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

See Also

[postgres](#)

pg_test_timing

pg_test_timing — measure timing overhead

Synopsis

pg_test_timing [*option...*]

Description

pg_test_timing is a tool to measure the timing overhead on your system and confirm that the system time never moves backwards. Systems that are slow to collect timing data can give less accurate EXPLAIN ANALYZE results.

Options

pg_test_timing accepts the following command-line options:

-d *duration*

--duration=*duration*

Specifies the test duration, in seconds. Longer durations give slightly better accuracy, and are more likely to discover problems with the system clock moving backwards. The default test duration is 3 seconds.

-V

--version

Print the pg_test_timing version and exit.

-?

--help

Show help about pg_test_timing command line arguments, and exit.

Usage

Interpreting Results

Good results will show most (>90%) individual timing calls take less than one microsecond. Average per loop overhead will be even lower, below 100 nanoseconds. This example from an Intel i7-860 system using a TSC clock source shows excellent performance:

Testing timing overhead for 3 seconds.

Per loop time including overhead: 35.96 ns

Histogram of timing durations:

< us	% of total	count
1	96.40465	80435604
2	3.59518	2999652
4	0.00015	126
8	0.00002	13
16	0.00000	2

Note that different units are used for the per loop time than the histogram. The loop can have resolution within a few nanoseconds (ns), while the individual timing calls can only resolve down to one microsecond (us).

Measuring Executor Timing Overhead

When the query executor is running a statement using EXPLAIN ANALYZE, individual operations are timed as well as showing a summary. The overhead of your system can be checked by counting rows with the psql program:

```
CREATE TABLE t AS SELECT * FROM generate_series(1,100000);
\timing
SELECT COUNT(*) FROM t;
EXPLAIN ANALYZE SELECT COUNT(*) FROM t;
```

The i7-860 system measured runs the count query in 9.8 ms while the `EXPLAIN ANALYZE` version takes 16.6 ms, each processing just over 100,000 rows. That 6.8 ms difference means the timing overhead per row is 68 ns, about twice what `pg_test_timing` estimated it would be. Even that relatively small amount of overhead is making the fully timed count statement take almost 70% longer. On more substantial queries, the timing overhead would be less problematic.

Changing Time Sources

On some newer Linux systems, it's possible to change the clock source used to collect timing data at any time. A second example shows the slowdown possible from switching to the slower `acpi_pm` time source, on the same system used for the fast results above:

```
# cat /sys/devices/system/clocksource/clocksource0/available_clocksource
tsc hpet acpi_pm
# echo acpi_pm > /sys/devices/system/clocksource/clocksource0/current_clocksource
# pg_test_timing
Per loop time including overhead: 722.92 ns
Histogram of timing durations:
  < us    % of total    count
    1      27.84870    1155682
    2      72.05956    2990371
    4       0.07810     3241
    8       0.01357      563
   16       0.00007       3
```

In this configuration, the sample `EXPLAIN ANALYZE` above takes 115.9 ms. That's 1061 ns of timing overhead, again a small multiple of what's measured directly by this utility. That much timing overhead means the actual query itself is only taking a tiny fraction of the accounted for time, most of it is being consumed in overhead instead. In this configuration, any `EXPLAIN ANALYZE` totals involving many timed operations would be inflated significantly by timing overhead.

FreeBSD also allows changing the time source on the fly, and it logs information about the timer selected during boot:

```
# dmesg | grep "Timecounter"
Timecounter "ACPI-fast" frequency 3579545 Hz quality 900
Timecounter "i8254" frequency 1193182 Hz quality 0
Timecounters tick every 10.000 msec
Timecounter "TSC" frequency 2531787134 Hz quality 800
# sysctl kern.timecounter.hardware=TSC
kern.timecounter.hardware: ACPI-fast -> TSC
```

Other systems may only allow setting the time source on boot. On older Linux systems the "clock" kernel setting is the only way to make this sort of change. And even on some more recent ones, the only option you'll see for a clock source is "jiffies". Jiffies are the older Linux software clock implementation, which can have good resolution when it's backed by fast enough timing hardware, as in this example:

```
$ cat /sys/devices/system/clocksource/clocksource0/available_clocksource
jiffies
$ dmesg | grep time.c
time.c: Using 3.579545 MHz WALL PM GTOD PIT/TSC timer.
time.c: Detected 2400.153 MHz processor.
$ pg_test_timing
Testing timing overhead for 3 seconds.
Per timing duration including loop overhead: 97.75 ns
Histogram of timing durations:
```

< us	% of total	count
1	90.23734	27694571
2	9.75277	2993204
4	0.00981	3010
8	0.00007	22
16	0.00000	1
32	0.00000	1

Clock Hardware and Timing Accuracy

Collecting accurate timing information is normally done on computers using hardware clocks with various levels of accuracy. With some hardware the operating systems can pass the system clock time almost directly to programs. A system clock can also be derived from a chip that simply provides timing interrupts, periodic ticks at some known time interval. In either case, operating system kernels provide a clock source that hides these details. But the accuracy of that clock source and how quickly it can return results varies based on the underlying hardware.

Inaccurate time keeping can result in system instability. Test any change to the clock source very carefully. Operating system defaults are sometimes made to favor reliability over best accuracy. And if you are using a virtual machine, look into the recommended time sources compatible with it. Virtual hardware faces additional difficulties when emulating timers, and there are often per operating system settings suggested by vendors.

The Time Stamp Counter (TSC) clock source is the most accurate one available on current generation CPUs. It's the preferred way to track the system time when it's supported by the operating system and the TSC clock is reliable. There are several ways that TSC can fail to provide an accurate timing source, making it unreliable. Older systems can have a TSC clock that varies based on the CPU temperature, making it unusable for timing. Trying to use TSC on some older multicore CPUs can give a reported time that's inconsistent among multiple cores. This can result in the time going backwards, a problem this program checks for. And even the newest systems can fail to provide accurate TSC timing with very aggressive power saving configurations.

Newer operating systems may check for the known TSC problems and switch to a slower, more stable clock source when they are seen. If your system supports TSC time but doesn't default to that, it may be disabled for a good reason. And some operating systems may not detect all the possible problems correctly, or will allow using TSC even in situations where it's known to be inaccurate.

The High Precision Event Timer (HPET) is the preferred timer on systems where it's available and TSC is not accurate. The timer chip itself is programmable to allow up to 100 nanosecond resolution, but you may not see that much accuracy in your system clock.

Advanced Configuration and Power Interface (ACPI) provides a Power Management (PM) Timer, which Linux refers to as the `acpi_pm`. The clock derived from `acpi_pm` will at best provide 300 nanosecond resolution.

Timers used on older PC hardware include the 8254 Programmable Interval Timer (PIT), the real-time clock (RTC), the Advanced Programmable Interrupt Controller (APIC) timer, and the Cyclone timer. These timers aim for millisecond resolution.

See Also

[EXPLAIN](#)

pg_waldump

`pg_waldump` — display a human-readable rendering of the write-ahead log of a Postgres Pro database cluster

Synopsis

```
pg_waldump [option...] [timestamp-option...] [startseg [endseg]]
```

Description

`pg_waldump` displays the write-ahead log (WAL) and prints timestamps for WAL records. This utility is mainly useful for debugging or educational purposes.

This utility can only be run by the user who installed the server, because it requires read-only access to the data directory.

Options

The following command-line options control the location and format of the output:

startseg

Start reading at the specified WAL segment file. This implicitly determines the path in which files will be searched for, and the timeline to use.

endseg

Stop after reading the specified WAL segment file.

`-b`

`--bkp-details`

Output detailed information about backup blocks.

`-B block`

`--block=block`

Only display records that modify the given block. The relation must also be provided with `--relation` or `-R`.

`-e end`

`--end=end`

Stop reading at the specified WAL location, instead of reading to the end of the log stream.

`-f`

`--follow`

After reaching the end of valid WAL, keep polling once per second for new WAL to appear.

`-F fork`

`--fork=fork`

Only display records that modify blocks in the given fork. The valid values are `main` for the main fork, `fsm` for the free space map, `vm` for the visibility map, and `init` for the init fork.

`-n limit`

`--limit=limit`

Display the specified number of records, then stop.

`-p path`
`--path=path`

Specifies a directory to search for WAL segment files or a directory with a `pg_wal` subdirectory that contains such files. The default is to search in the current directory, the `pg_wal` subdirectory of the current directory, and the `pg_wal` subdirectory of `PGDATA`.

`-q`
`--quiet`

Do not print any output, except for errors. This option can be useful when you want to know whether a range of WAL records can be successfully parsed but don't care about the record contents.

`-r rmgr`
`--rmgr=rmgr`

Only display records generated by the specified resource manager. You can specify the option multiple times to select multiple resource managers. If `list` is passed as name, print a list of valid resource manager names, and exit.

Extensions may define custom resource managers, but `pg_waldump` does not load the extension module and therefore does not recognize custom resource managers by name. Instead, you can specify the custom resource managers as `custom###` where `###` is the three-digit resource manager ID. Names of this form will always be considered valid.

`-R tblspc/db/rel`
`--relation=tblspc/db/rel`

Only display records that modify blocks in the given relation. The relation is specified with tablespace OID, database OID, and relfilenode separated by slashes, for example `1234/12345/12345`. This is the same format used for relations in the program's output.

`-s start`
`--start=start`

WAL location at which to start reading. The default is to start reading the first valid WAL record found in the earliest file found.

`-t timeline`
`--timeline=timeline`

Timeline from which to read WAL records. The default is to use the value in `startseg`, if that is specified; otherwise, the default is 1. The value can be specified in decimal or hexadecimal, for example `17` or `0x11`.

`-V`
`--version`

Print the `pg_waldump` version and exit.

`-w`
`--fullpage`

Only display records that include full page images.

`-x xid`
`--xid=xid`

Only display records marked with the given transaction ID.

`-z`
`--stats[=record]`

Display summary statistics (number and size of records and full-page images) instead of individual records. Optionally generate statistics per-record instead of per-rmgr.

If `pg_waldump` is terminated by signal SIGINT (**Control+C**), the summary of the statistics computed is displayed up to the termination point. This operation is not supported on Windows.

`--save-fullpage=save_path`

Save full page images found in the WAL records to the `save_path` directory. The images saved are subject to the same filtering and limiting criteria as the records displayed.

The full page images are saved with the following file name format: `TIMELINE-LSN.RELTABLESPACE.DATOID.RELNODE.BLKNO_FORK`. The file names are composed of the following parts:

Component	Description
TIMELINE	The timeline of the WAL segment file where the record is located formatted as one 8-character hexadecimal number <code>%08X</code>
LSN	The LSN of the record with this image, formatted as two 8-character hexadecimal numbers <code>%08X-%08X</code>
RELTABLESPACE	tablespace OID of the block
DATOID	database OID of the block
RELNODE	filenode of the block
BLKNO	block number of the block
FORK	The name of the fork the full page image came from, such as <code>main</code> , <code>fsm</code> , <code>vm</code> , or <code>init</code> .

`-?`

`--help`

Show help about `pg_waldump` command line arguments, and exit.

The following command-line options enable printing timestamps for various types of WAL records. You can use these options together with `startseg`, `endseg`, `-s`, and `-e` options.

`-E`

`--end-timestamp`

Print the timestamp of the last WAL record of the specified type found in the log segment file. When using this option, you must also specify the `-s` option.

By default, `pg_waldump` prints timestamps only for COMMIT records. You can specify other record types using the `--timestamp-filter` option.

`-S`

`--start-timestamp`

Print the timestamp of the first WAL record of the specified type found in the log segment file. This option is required if you are going to use `-E` or `--timestamp-filter` options.

By default, `pg_waldump` prints timestamps only for COMMIT records. You can specify other record types using the `--timestamp-filter` option.

`--timestamp-filter=argument [, ...]`

Specify WAL record types for which to print timestamps. When using this option, you must also specify the `-s` option.

The `--timestamp-filter` option can take the following arguments, in the comma-separated format:

- `XLOG_RESTORE_POINT` — named restore points created with the `pg_create_restore_point()` function.

- `XLOG_XACT_COMMIT` — commit records for transactions. These records are caused by the [COMMIT](#) command.
- `XLOG_XACT_COMMIT_PREPARED` — commit records for transactions that were earlier prepared for a two-phase commit. These records are caused by the [COMMIT PREPARED](#) command.
- `XLOG_XACT_ABORT` — abort records for transactions. These records are caused by the [ROLLBACK](#) command.
- `XLOG_XACT_ABORT_PREPARED` — abort records for transactions that were earlier prepared for a two-phase commit. These records are caused by the [ROLLBACK PREPARED](#) command.

By default, `pg_waldump` prints timestamps only for COMMIT records.

Environment

`PGDATA`

Data directory; see also the `-p` option.

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

Can give wrong results when the server is running.

Only the specified timeline is displayed (or the default, if none is specified). Records in other timelines are ignored.

`pg_waldump` cannot read WAL files with suffix `.partial`. If those files need to be read, `.partial` suffix needs to be removed from the file name.

See Also

[Section 28.7](#)

pg_walsummary

pg_walsummary — print contents of WAL summary files

Synopsis

```
pg_walsummary [option...] [file...]
```

Description

pg_walsummary is used to print the contents of WAL summary files. These binary files are found with the `pg_wal/summaries` subdirectory of the data directory, and can be converted to text using this tool. This is not ordinarily necessary, since WAL summary files primarily exist to support incremental backup, but it may be useful for debugging purposes.

A WAL summary file is indexed by tablespace OID, relation OID, and relation fork. For each relation fork, it stores the list of blocks that were modified by WAL within the range summarized in the file. It can also store a "limit block," which is 0 if the relation fork was created or truncated within the relevant WAL range, and otherwise the shortest length to which the relation fork was truncated. If the relation fork was not created, deleted, or truncated within the relevant WAL range, the limit block is undefined or infinite and will not be printed by this tool.

Options

```
-i  
--individual
```

By default, `pg_walsummary` prints one line of output for each range of one or more consecutive modified blocks. This can make the output a lot briefer, since a relation where all blocks from 0 through 999 were modified will produce only one line of output rather than 1000 separate lines. This option requests a separate line of output for every modified block.

```
-q  
--quiet
```

Do not print any output, except for errors. This can be useful when you want to know whether a WAL summary file can be successfully parsed but don't care about the contents.

```
-V  
--version
```

Display version information, then exit.

```
-?  
--help
```

Shows help about `pg_walsummary` command line arguments, and exits.

Environment

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

See Also

[pg_combinebackup](#)

postgres

postgres — Postgres Pro database server

Synopsis

postgres [*option...*]

Description

postgres is the Postgres Pro database server. In order for a client application to access a database it connects (over a network or locally) to a running postgres instance. The postgres instance then starts a separate server process to handle the connection.

One postgres instance always manages the data of exactly one database cluster. A database cluster is a collection of databases that is stored at a common file system location (the “data area”). More than one postgres instance can run on a system at one time, so long as they use different data areas and different communication ports (see below). When postgres starts it needs to know the location of the data area. The location must be specified by the `-D` option or the `PGDATA` environment variable; there is no default. Typically, `-D` or `PGDATA` points directly to the data area directory created by [initdb](#). Other possible file layouts are discussed in [Section 19.2](#).

By default postgres starts in the foreground and prints log messages to the standard error stream. In practical applications postgres should be started as a background process, perhaps at boot time.

The postgres command can also be called in single-user mode. The primary use for this mode is during bootstrapping by [initdb](#). Sometimes it is used for debugging or disaster recovery; note that running a single-user server is not truly suitable for debugging the server, since no realistic interprocess communication and locking will happen. When invoked in single-user mode from the shell, the user can enter queries and the results will be printed to the screen, but in a form that is more useful for developers than end users. In the single-user mode, the session user will be set to the user with ID 1, and implicit superuser powers are granted to this user. This user does not actually have to exist, so the single-user mode can be used to manually recover from certain kinds of accidental damage to the system catalogs.

Options

postgres accepts the following command-line arguments. For a detailed discussion of the options consult [Chapter 19](#). You can save typing most of these options by setting up a configuration file. Some (safe) options can also be set from the connecting client in an application-dependent way to apply only for that session. For example, if the environment variable `PGOPTIONS` is set, then libpq-based clients will pass that string to the server, which will interpret it as postgres command-line options.

General Purpose

`-B nbuffers`

Sets the number of shared buffers for use by the server processes. The default value of this parameter is chosen automatically by [initdb](#). Specifying this option is equivalent to setting the [shared_buffers](#) configuration parameter.

`-c name=value`

Sets a named run-time parameter. The configuration parameters supported by Postgres Pro are described in [Chapter 19](#). Most of the other command line options are in fact short forms of such a parameter assignment. `-c` can appear multiple times to set multiple parameters.

`-C name`

Prints the value of the named run-time parameter, and exits. (See the `-c` option above for details.) This returns values from `postgresql.conf`, modified by any parameters supplied in this invocation. It does not reflect parameters supplied when the cluster was started.

This can be used on a running server for most parameters. However, the server must be shut down for some runtime-computed parameters (e.g., [shared_memory_size](#), [shared_memory_size_in_huge_pages](#), and [wal_segment_size](#)).

This option is meant for other programs that interact with a server instance, such as [pg_ctl](#), to query configuration parameter values. User-facing applications should instead use [SHOW](#) or the `pg_settings` view.

`-d debug-level`

Sets the debug level. The higher this value is set, the more debugging output is written to the server log. Values are from 1 to 5. It is also possible to pass `-d 0` for a specific session, which will prevent the server log level of the parent `postgres` process from being propagated to this session.

`-D datadir`

Specifies the file system location of the database configuration files. See [Section 19.2](#) for details.

`-e`

Sets the default date style to “European”, that is `DMY` ordering of input date fields. This also causes the day to be printed before the month in certain date output formats. See [Section 8.5](#) for more information.

`-F`

Disables `fsync` calls for improved performance, at the risk of data corruption in the event of a system crash. Specifying this option is equivalent to disabling the [fsync](#) configuration parameter. Read the detailed documentation before using this!

`-h hostname`

Specifies the IP host name or address on which `postgres` is to listen for TCP/IP connections from client applications. The value can also be a comma-separated list of addresses, or `*` to specify listening on all available interfaces. An empty value specifies not listening on any IP addresses, in which case only Unix-domain sockets can be used to connect to the server. Defaults to listening only on localhost. Specifying this option is equivalent to setting the [listen_addresses](#) configuration parameter.

`-i`

Allows remote clients to connect via TCP/IP (Internet domain) connections. Without this option, only local connections are accepted. This option is equivalent to setting `listen_addresses` to `*` in `postgresql.conf` or via `-h`.

This option is deprecated since it does not allow access to the full functionality of [listen_addresses](#). It's usually better to set `listen_addresses` directly.

`-k directory`

Specifies the directory of the Unix-domain socket on which `postgres` is to listen for connections from client applications. The value can also be a comma-separated list of directories. An empty value specifies not listening on any Unix-domain sockets, in which case only TCP/IP sockets can be used to connect to the server. The default value is normally `/tmp`, but that can be changed at build time. Specifying this option is equivalent to setting the [unix_socket_directories](#) configuration parameter.

`-l`

Enables secure connections using SSL. Postgres Pro must have been compiled with support for SSL for this option to be available. For more information on using SSL, refer to [Section 18.9](#).

`-N max-connections`

Sets the maximum number of client connections that this server will accept. The default value of this parameter is chosen automatically by `initdb`. Specifying this option is equivalent to setting the [max_connections](#) configuration parameter.

`-p port`

Specifies the TCP/IP port or local Unix domain socket file extension on which `postgres` is to listen for connections from client applications. Defaults to the value of the `PGPORT` environment variable, or if `PGPORT` is not set, then defaults to the value established during compilation (normally 5432). If you specify a port other than the default port, then all client applications must specify the same port using either command-line options or `PGPORT`.

`-s`

Print time information and other statistics at the end of each command. This is useful for benchmarking or for use in tuning the number of buffers.

`-S work-mem`

Specifies the base amount of memory to be used by sorts and hash tables before resorting to temporary disk files. See the description of the `work_mem` configuration parameter in [Section 19.4.1](#).

`-V`

`--version`

Print the postgres version and exit.

`--name=value`

Sets a named run-time parameter; a shorter form of `-c`.

`--describe-config`

This option dumps out the server's internal configuration variables, descriptions, and defaults in tab-delimited `COPY` format. It is designed primarily for use by administration tools.

`-?`

`--help`

Show help about postgres command line arguments, and exit.

Semi-Internal Options

The options described here are used mainly for debugging purposes, and in some cases to assist with recovery of severely damaged databases. There should be no reason to use them in a production database setup. They are listed here only for use by Postgres Pro system developers. Furthermore, these options might change or be removed in a future release without notice.

`-f { s | i | o | b | t | n | m | h }`

Forbids the use of particular scan and join methods: `s` and `i` disable sequential and index scans respectively, `o`, `b` and `t` disable index-only scans, bitmap index scans, and TID scans respectively, while `n`, `m`, and `h` disable nested-loop, merge and hash joins respectively.

Neither sequential scans nor nested-loop joins can be disabled completely; the `-fs` and `-fn` options simply discourage the optimizer from using those plan types if it has any other alternative.

`-O`

Allows the structure of system tables to be modified. This is used by `initdb`.

`-P`

Ignore system indexes when reading system tables, but still update the indexes when modifying the tables. This is useful when recovering from damaged system indexes.

`-t pa[rser] | pl[anner] | e[xecutor]`

Print timing statistics for each query relating to each of the major system modules. This option cannot be used together with the `-s` option.

-T

This option is for debugging problems that cause a server process to die abnormally. The ordinary strategy in this situation is to notify all other server processes that they must terminate, by sending them SIGQUIT signals. With this option, SIGABRT will be sent instead, resulting in production of core dump files.

-v *protocol*

Specifies the version number of the frontend/backend protocol to be used for a particular session. This option is for internal use only.

-W *seconds*

A delay of this many seconds occurs when a new server process is started, after it conducts the authentication procedure. This is intended to give an opportunity to attach to the server process with a debugger.

-Z

Verifies that the current postgres binary is compatible with the specified cluster. If the architecture type or any compilation options that affect cluster compatibility do not match, returns an exit status of 1 and provides an error message specifying the first detected incompatibility. Otherwise, reports success and exits without starting the cluster.

You must provide the path to the data directory of the cluster to check.

The cluster and the binary must have the same byte order and architecture type for this option to work correctly.

Tip

If the binary appears to be incompatible with the specified cluster, run [pgpro_controldata](#) with the `-P` or `-C` command-line argument to view all parameters of the cluster that affect the compatibility.

Options for Single-User Mode

The following options only apply to the single-user mode (see [Single-User Mode](#) below).

--single

Selects the single-user mode. This must be the first argument on the command line.

database

Specifies the name of the database to be accessed. This must be the last argument on the command line. If it is omitted it defaults to the user name.

-E

Echo all commands to standard output before executing them.

-j

Use semicolon followed by two newlines, rather than just newline, as the command entry terminator.

-r *filename*

Send all server log output to *filename*. This option is only honored when supplied as a command-line option.

Environment

PGCLIENTENCODING

Default character encoding used by clients. (The clients can override this individually.) This value can also be set in the configuration file.

PGDATA

Default data directory location

PGDATESTYLE

Default value of the [DateStyle](#) run-time parameter. (The use of this environment variable is deprecated.)

PGPORT

Default port number (preferably set in the configuration file)

Diagnostics

A failure message mentioning `semget` or `shmget` probably indicates you need to configure your kernel to provide adequate shared memory and semaphores. For more discussion see [Section 18.4](#). You might be able to postpone reconfiguring your kernel by decreasing [shared_buffers](#) to reduce the shared memory consumption of Postgres Pro, and/or by reducing [max_connections](#) to reduce the semaphore consumption.

A failure message suggesting that another server is already running should be checked carefully, for example by using the command

```
$ ps ax | grep postgres
```

or

```
$ ps -ef | grep postgres
```

depending on your system. If you are certain that no conflicting server is running, you can remove the lock file mentioned in the message and try again.

A failure message indicating inability to bind to a port might indicate that that port is already in use by some non-Postgres Pro process. You might also get this error if you terminate `postgres` and immediately restart it using the same port; in this case, you must simply wait a few seconds until the operating system closes the port before trying again. Finally, you might get this error if you specify a port number that your operating system considers to be reserved. For example, many versions of Unix consider port numbers under 1024 to be “trusted” and only permit the Unix superuser to access them.

Notes

The utility command [pg_ctl](#) can be used to start and shut down the `postgres` server safely and comfortably.

If at all possible, *do not* use `SIGKILL` to kill the main `postgres` server. Doing so will prevent `postgres` from freeing the system resources (e.g., shared memory and semaphores) that it holds before terminating. This might cause problems for starting a fresh `postgres` run.

To terminate the `postgres` server normally, the signals `SIGTERM`, `SIGINT`, or `SIGQUIT` can be used. The first will wait for all clients to terminate before quitting, the second will forcefully disconnect all clients, and the third will quit immediately without proper shutdown, resulting in a recovery run during restart.

The `SIGHUP` signal will reload the server configuration files. It is also possible to send `SIGHUP` to an individual server process, but that is usually not sensible.

To cancel a running query, send the `SIGINT` signal to the process running that command. To terminate a backend process cleanly, send `SIGTERM` to that process. See also `pg_cancel_backend` and `pg_terminate_backend` in [Section 9.28.2](#) for the SQL-callable equivalents of these two actions.

The `postgres` server uses `SIGQUIT` to tell subordinate server processes to terminate without normal cleanup. This signal *should not* be used by users. It is also unwise to send `SIGKILL` to a server process — the main `postgres` process will interpret this as a crash and will force all the sibling processes to quit as part of its standard crash-recovery procedure.

Bugs

The `--` options will not work on FreeBSD or OpenBSD. Use `-c` instead. This is a bug in the affected operating systems; a future release of Postgres Pro will provide a workaround if this is not fixed.

Single-User Mode

To start a single-user mode server, use a command like

```
postgres --single -D /usr/local/pgsql/data other-options my_database
```

Provide the correct path to the database directory with `-D`, or make sure that the environment variable `PGDATA` is set. Also specify the name of the particular database you want to work in.

Normally, the single-user mode server treats newline as the command entry terminator; there is no intelligence about semicolons, as there is in `psql`. To continue a command across multiple lines, you must type backslash just before each newline except the last one. The backslash and adjacent newline are both dropped from the input command. Note that this will happen even when within a string literal or comment.

But if you use the `-j` command line switch, a single newline does not terminate command entry; instead, the sequence semicolon-newline-newline does. That is, type a semicolon immediately followed by a completely empty line. Backslash-newline is not treated specially in this mode. Again, there is no intelligence about such a sequence appearing within a string literal or comment.

In either input mode, if you type a semicolon that is not just before or part of a command entry terminator, it is considered a command separator. When you do type a command entry terminator, the multiple statements you've entered will be executed as a single transaction.

To quit the session, type EOF (**Control+D**, usually). If you've entered any text since the last command entry terminator, then EOF will be taken as a command entry terminator, and another EOF will be needed to exit.

Note that the single-user mode server does not provide sophisticated line-editing features (no command history, for example). Single-user mode also does not do any background processing, such as automatic checkpoints or replication.

Examples

To start `postgres` in the background using default values, type:

```
$ nohup postgres >logfile 2>&1 </dev/null &
```

To start `postgres` with a specific port, e.g., 1234:

```
$ postgres -p 1234
```

To connect to this server using `psql`, specify this port with the `-p` option:

```
$ psql -p 1234
```

or set the environment variable `PGPORT`:

```
$ export PGPORT=1234
$ psql
```

Named run-time parameters can be set in either of these styles:

```
$ postgres -c work_mem=1234  
$ postgres --work-mem=1234
```

Either form overrides whatever setting might exist for `work_mem` in `postgresql.conf`. Notice that underscores in parameter names can be written as either underscore or dash on the command line. Except for short-term experiments, it's probably better practice to edit the setting in `postgresql.conf` than to rely on a command-line switch to set a parameter.

See Also

[initdb](#), [pg_ctl](#)

shardmanctl

shardmanctl — Postgres Pro Shardman auxiliary command-line client and deployment tool

Table of Contents

[Description](#)

[Backup and Restore](#)

[Configuration-related Commands](#)

[Cluster-related Commands](#)

[Nodes-related Commands](#)

[Shard-related Commands](#)

[Status-related Commands](#)

[Store-related Commands](#)

[Table-related Commands](#)

[Bench-related Commands](#)

[Cleanup-related Commands](#)

[Schema-related Commands](#)

[Daemon-related Commands](#)

[Other Commands](#)

[Common Options](#)

[Environment](#)

[Examples](#)

Description

shardmanctl is an utility for managing a Shardman cluster.

For any command that uses the node name as an argument, the node name can be specified either by its hostname or IP address.

The `backup` command is used to backup a Postgres Pro Shardman cluster. A backup consists of a directory with base backups of all replication groups and WAL files needed for recovery. etcd metadata is saved to the `etcd_dump` file. The `backup_info` file is created during a backup and contains the backup description.

The `cleanup` command is used for cleanup after failure of the `nodes add` command or of the [shardmanctl rebalance](#) command. Final changes to the etcd store are done at the end of the command execution. This simplifies the `cleanup` process. During cleanup, incomplete clover definitions and definitions of the corresponding replication groups are removed from the etcd metadata. Definitions of the corresponding foreign servers are removed from the DBMS metadata of the remaining replication groups. Since the `cleanup` process can be destructive, by default, the tool operates in the report-only mode: it only shows actions to be done during the actual cleanup. To perform the actual cleanup, add the `-p` flag. For usage details of the command, see [the section called “Performing Cleanup”](#).

The `daemon check` command is used to verify that shardmand daemon is running on the nodes specified by `--nodes` option and is configured for the same cluster as `shardmanctl`. For usage details of the command, see [the section called “Checking shardmand Service on Nodes”](#).

The `init` command is used to register a new Postgres Pro Shardman cluster in the etcd store or to reinitialize the existing cluster defining a new cluster configuration and removing all data and nodes. In the `init` mode, `shardmanctl` reads the cluster specification, processes it and saves to the etcd store as parts of two JSON documents: `ClusterSpec` — as part of `shardman/cluster0/data/cluster` and `LadleSpec` — as part of `shardman/cluster0/data/ladle` (`cluster0` is the default cluster name used by Postgres Pro Shardman utilities). Common options related to the etcd store, such as `--store-endpoints`, are also saved to the etcd store and pushed down to all Postgres Pro Shardman services started by `shardmand`. For the description of the Postgres Pro Shardman initialization file format, see [Section 19.21.2](#). For usage details of the command, see [the section called “Registering a Postgres Pro Shardman Cluster”](#).

The `config generate` command is used to create a default `sdmspec.json` template. By default, data is returned to the standard output. To write the result to a file, use flag `-f filename`. For the description of the Shardman initialization file format, see [Section 19.21.2](#).

The `config verify` command is used to check a correctness of the input Postgres Pro Shardman initialization file. By default, the configuration is read from standard input. To read the configuration from a file, use flag `-f filename`. For the description of the Postgres Pro Shardman initialization file format, see [Section 19.21.2](#).

The `config get` command is used to output the current full cluster specification or a configuration of the specified revision. The command takes the current cluster configuration from the cluster store. For the description of the Postgres Pro Shardman initialization file format, see [Section 19.21.2](#).

The `config update` command is used to update the BiHA cluster or full Postgres Pro Shardman configuration. The new configuration is applied to all replication groups and is saved in `shardman/cluster0/data/cluster` etcd key. Note that `config update` can cause a DBMS restart.

The `forall` command is used to execute an SQL statement on all replication groups in a Postgres Pro Shardman cluster.

The `getconnstr` command is used to get the libpq connection string for connecting to a cluster as administrator.

The `load` command is used to upload data from a text file to a distributed table or to upload a database schema from a PostgreSQL database to Postgres Pro Shardman. When loading data from a file, `text` and `csv` formats are supported. If a file is compressed with `gzip`, it will be automatically decoded while reading. To read data from `stdin`, specify `--file=-`. The data loading process can be optimized by specifying the number of parallel workers (`key -j`).

The `nodes add` command is used to add new nodes to a Shardman cluster. With the default `cross` placement policy, nodes are added to a cluster by `clovers`. Each node in a clover runs the primary DBMS instance and perhaps several replicas of other nodes in the clover. The number of replicas is determined by the [Repfactor](#) configuration parameter. So, each clover consists of `Repfactor + 1` nodes and can stand loss of `Repfactor` nodes.

With `manual` placement policy, each new node is added as a replication group consisting of one primary server. After adding primary nodes, you can add replicas to the new replication group by calling the `shard add` command.

`shardmanctl` performs the `nodes add` operation in several steps:

1. Acquires a global metadata lock.
2. For each specified node, checks that Postgres Pro Shardman is running on it and that it sees the current cluster configuration.

3. Calculates the services to be present on each node and saves this information in etcd as part of the `shardman/cluster0/data/ladle Layout` object.
4. Generates the configuration for new BiHA clusters (also called *replication groups*) and initializes them.
5. Registers the added replication groups in the `shardman/cluster0/data/ladle etcd` key.
6. Waits for shardmand to start all the necessary services, checks that new replication groups are accessible and have correct configuration.
7. Creates an auxiliary broadcaster that holds locks on each existing replication group in the cluster.
8. For each new replication group, copies all schemas and `shardman` schema data from a randomly selected existing replication group to the new one, ensures that the Postgres Pro Shardman extension is installed on the new replication group, and recalculates OIDs used in the extension configuration tables.
9. On each existing replication group, defines foreign servers referencing the new replication group and recreates definitions of foreign servers on the new replication group.
10. Recreates all partitions of sharded tables as foreign tables referencing data from old replication groups and has the changes registered in the etcd storage.
11. For each new replication group, copies the global table data from existing replication groups to the new one.
12. Rebalances partitions of sharded tables. The rebalancing process for each sharded table iteratively determines the replication group with the maximum and minimum number of partitions and creates a task to move one partition to the replication group with the minimum number of partitions. This process is repeated while `max - min > 1`. To move partitions, we use logical replication. Partitions of colocated tables are moved together with partitions of the distributed tables to which they refer. You can skip this step using the `--no-rebalance`.

For usage details of the command, see [the section called “Adding Nodes to a Postgres Pro Shardman Cluster”](#).

The `nodes rm` command is used to remove nodes from a Shardman cluster. In the manual-topology mode, this command only removes the specified nodes from the cluster and if a node is the last in the replication group, the entire group gets removed. In the cross-replication mode, this command removes clovers containing the specified nodes from the cluster. The last clover in the cluster cannot be removed. Any data (such as partitions of sharded relations) on removed replication groups is migrated to the remaining replication groups using logical replication, and all references to the removed replication groups (including definitions of foreign servers) are removed from the metadata of the remaining replication groups. Finally, the metadata in etcd is updated. For usage details of the command, see [the section called “Removing Nodes from a Postgres Pro Shardman cluster”](#).

The `probackup` command is used to backup and restore the Shardman cluster using `pg_probackup` backup utility. For usage details of the command, see [the section called “probackup Command”](#).

The `rebalance` command is used to evenly rebalance sharded tables in a cluster. This can be useful, for example, if you did not perform rebalance when adding nodes to the cluster. If the `--force` option is not provided, then tables with manually moved partitions will be skipped.

The `cleanup` command with flag `--after-rebalance` is used to perform cleanup after failure of a `rebalance` command. On each node, it cleans up subscriptions and publications left from the `rebalance` command and drops tables that store data of partially-transferred partitions of sharded tables.

The `cluster repfactor set` command is used to set the value of the replication factor for the Postgres Pro Shardman cluster. This command can only be used in `manual` topology cluster mode. The value of the new replication factor is passed through the command line flag `--value repfactor`.

The `cluster start` command is used to start all stopped PostgreSQL instances with the `cluster stop` command. For the command to work, `shardmand` must be running.

The `cluster stop` command is used to stop all PostgreSQL instances for the Postgres Pro Shardman cluster. At the same time, the shardmand daemons continue to work.

The `cluster topology` command is used to visualize the topology of a cluster. By default, the topology is returned in a table view. If you want to get a JSON or text representation, then use the flag `--format json|text`.

The `recover` command is used to restore a Postgres Pro Shardman cluster from a backup created by the `backup` command.

The `restart` command is used to restart a Postgres Pro Shardman cluster, including all shardmand instances. If PostgreSQL instances were previously stopped using the `cluster stop` command, they will be started. The command returns control after all primary nodes in the cluster have been restarted.

The `set` command is used to set one or more parameters for DBMS instances of the Postgres Pro Shardman cluster. Parameters are passed as arguments to the command line, each of them looks like `param=value`. The command is actually an alternative to `shardmanctl config update -p` to update database settings.

The `status` command is used to display health status of Postgres Pro Shardman cluster subsystems. It can show status of several components: store, metadata, shardmand, replication groups, primary nodes, dictionary, status from `biha.status_v`, and restart of the required parameters. If only some subsystems are of interest, option `--filter` may be used. Also `status` supports sorting its messages by `status`, `node` or `replication group` and printing the result to stdout as a table (`table`), text (`text`) or JSON (`json`) with `table` as the default. For usage details of the command, see [the section called “Getting the Status of Cluster Subsystems”](#).

The `store dump` command gets all the keys and their values from the etcd store and outputs them into the `--file`, where `-` value is used for outputting to stdout (default). It is intended to be used for debugging, so some harmless errors may be produced during execution, yet all the available information will be dumped. Only keys for the current cluster (with current cluster prefix like `shardman/cluster0/`) will be dumped. For usage details of the command, see [the section called “Dumping All Keys from the Store to Debug Error Configuration”](#).

The `store get` command gets a particular value from the store by its key name. It is expected to be a JSON value, so if it is not (which is not prohibited), some harmless errors may be produced. The key to retrieve from store can be specified with `--key` option; several keys have *aliases* — short names for easy use. To get a key by its alias, use `--alias` option with one of the available aliases (use `--help` or examples below for reference). Also aliases `shardspec` and `spec` can be used to manipulate initial cluster and BiHA configuration explicitly, without retrieving it from the full cluster specification. It is recommended to use existing aliases instead of full key names since there are some additional checks in alias processing, which help to achieve safer results. By default, a key is printed to stdout (explicitly — with `--file=-` option), but can be output to any desired file. For usage details of the command, see [the section called “Getting the Current Shard Specification”](#).

The `store keys` command shows all the keys in the store for the current cluster (with cluster prefix) and its aliases. Aliases `shardspec` and `spec` are not shown since they are parts of other keys. For usage details of the command, see [the section called “Getting the Cluster and Ladle Key Names For the Current Cluster”](#).

The `store set` command creates or rewrites one particular key in the store. It is not expected to be a JSON value for a random key, but if it is one of the keys that have aliases with a known mapping (like `ladle` or `cluster`), the command will not accept incorrect JSON structures. Just like `store get` command, `store set` accepts a key name via `--key` or `--alias` option and the input source file as `--file` (stdin is specified with `-` value). For usage details of the command, see [the section called “Setting a New Spec for the Cluster”](#).

The `store lock` command shows the current cluster meta lock information. In case lock does not exist returns `Lock not found`. Displays cluster id, command that acquired locks, host name and lock time.

You can specify `--format` to output in `json` format or in `text` format (by default). For usage details of the command, see [the section called “Output Current Cluster Meta Lock Information”](#).

The `upgrade` command is used to update the version of Postgresql shardman extension on all cluster nodes. Before upgrading extensions, you need to install new packages and run the `restart` command. As a result of `upgrade`, utilities will upgrade shardman and all the other extensions on the server.

Sometimes after running the `upgrade` command or some user's manual manipulations, dictionary errors may appear in the output of the `status` command. One of the reasons for these errors is that the value of the `srvoptions` field of the `pg_foreign_server` table differs from what the system expects. To solve this specific issue, use the `config update fdw` command, which will return `srvoptions` to the expected state.

Note

Most of the described `shardmanctl` commands take a global metadata lock.

Backup and Restore-related Commands

backup Command

To backup a Postgres Pro Shardman cluster, you can run the following command:

```
shardmanctl [common_options] backup --datadir directory [--use-ssh]
```

You must pass the directory to write the output to through the `--datadir` option. You can limit the number of running concurrent tasks (`pg_receivewal` or `pg_basebackup` commands) by passing the limit through the `--maxtasks` option.

If `--use-ssh` is specified [shardmanctl recover](#) command will use `scp` command to restore data. It allows to use backup repository on the local host.

probackup Command

Syntax:

```
shardmanctl [common_options] probackup
    [init|archive-command|backup|checkdb|delete|merge|restore|set-config|show|
validate|show-config]
    [--log-to-console] [--help]
    [subcommand_options]
```

Creates a backup of a Shardman cluster and restores the Postgres Pro Shardman cluster from a backup using `pg_probackup`.

List of subcommands:

`init`

Initializes a new repository folder for the Postgres Pro Shardman cluster backup and creates a configuration file on all nodes for connection to the backup storage if `--storage-type` is `S3`.

`archive-command`

Adds `archive_command` to each replication group (or to a single one if the `--shard` option is specified) and enables or disables it in the Postgres Pro Shardman cluster.

`backup`

Creates a backup of the Postgres Pro Shardman cluster.

checkdb

Verifies the Postgres Pro Shardman cluster correctness by detecting physical and logical corruption.

delete

Deletes a backup of the Postgres Pro Shardman cluster with the specified `backup_id`.

merge

Merges the backups that belong to a common incremental backup chain. The full backup merges the backups with their first incremental backup. The incremental backup merges the backups with their parent full backup, along with all the incremental backups between them. Once the merge is complete, the full backup covers all the merged data, and the incremental backups are removed as redundant. In this version, you cannot run the `merge` command using the S3 interface.

restore

Restores the Postgres Pro Shardman cluster from the selected backup.

show

Shows the list of backups of the Postgres Pro Shardman cluster.

validate

Checks the selected Postgres Pro Shardman cluster backup for integrity.

show-config

Displays all the current `pg_probackup` configuration settings, including those that are specified in the `pg_probackup.conf` configuration file located in the `backup_dir/backups/shard_name` directory and those that were provided on a command line.

set-config

Adds the specified settings to the `pg_probackup.conf` or modifies those previously added.

The following options can be used with all `probackup` subcommands:

--log-to-console

Outputs a full `probackup` log to the console. By default, for each replication group the `probackup` log file is written to the backup directory (see `--backup-path` below) as the `<backup-directory>/backup/log/pg_probackup-<repgroup-name>.log` file. The log rotation file size is 20MB. If this value is reached, the log file is rotated once a `shardmanctl probackup validate` or `shardmanctl probackup backup` command is launched.

--help

Shows subcommand help.

init**Syntax:**

```
shardmanctl [common_options] probackup init
-B|--backup-path path
-E|--etcd-path path
[--remote-port port]
[--remote-user username]
[--ssh-key path]
[-t|--timeout seconds]
[-m|--maxtasks number_of_tasks]
[--storage-type mount|remote|S3]
[--s3-config-only]
```

```
[--s3-config-path path]  
[--s3-host S3_host]  
[--s3-port S3_port]  
[--s3-access-key S3_access_key]  
[--s3-secret-key S3_secret_key]  
[--s3-bucket S3_bucket]  
[--s3-region S3_region]  
[--s3-buffer-size size]  
[--s3-retries number_of_retries]  
[--s3-timeout time]  
[--s3-https]  
[-y|--yes]
```

Initializes a new repository folder for the Postgres Pro Shardman cluster backup.

```
-B path  
--backup-path path
```

Required if `--s3-config-only` is not used. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

```
-E path  
--etcd-path path
```

Required if `--s3-config-only` is not used. Specifies the path to the catalog where the etcd dumps should be stored.

```
--remote-port port
```

Specifies the remote ssh port for replication group instances.

Default: 22.

```
--remote-user username
```

Specifies the remote ssh user for replication group instances.

Default: postgres.

```
--ssh-key path
```

Specifies the ssh private key for execution of remote ssh commands.

Default: `$HOME/.ssh/id_rsa`.

```
--storage-type mount|remote|S3
```

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively.

Default: `remote`.

```
--s3-config-path path
```

Specifies the path where the S3 configuration file will be created on all Postgres Pro Shardman nodes.

Default: `<shardman-data-dir>/s3.config`.

```
--s3-config-only
```

Create only S3 configuration files on all nodes and skip backup repository initialization. This flag is useful if the value of `--storage-type` is `S3`.

`--s3-host host`

Specifies the S3 host to connect to S3-compatible storage.

`--s3-port port`

Specifies the S3 port to connect to S3-compatible storage.

`--s3-access-key access-key`

Specifies the S3 access key to connect to S3-compatible storage.

`--s3-secret-key access-key`

Specifies the S3 secret key to connect to the S3-compatible storage.

`--s3-bucket bucket`

Specifies the bucket in the S3-compatible object storage for storing backups.

`--s3-region bucket`

Specifies the region in the S3-compatible object storage.

`--s3-buffer-size size`

Size of the read/write buffer for `pg_probackup` to communicate with the S3-compatible object storage, in MiB.

Default: 16.

`--s3-retries number_of_retries`

Maximum number of attempts for `pg_probackup` to execute an S3 request in case of failures.

Default: 5.

`--s3-timeout time`

Maximum allowable amount of time for `pg_probackup` to transfer data of size `--s3-buffer-size` to/from the S3-compatible object storage, in seconds.

Default: 300.

`--s3-https`

Specifies the HTTPS URL to connect to the S3-compatible object storage.

`-y|--yes`

Approve the operation regardless of whether the file specified in `--s3-config-path` exists.

archive-command

Syntax:

```
shardmanctl [common_options] probackup archive-command [add|rm]  
  -B|--backup-path path  
  [-j|--jobs count]  
  [--compress]  
  [--compress-algorithm algorithm]  
  [--compress-level level]  
  [--batch-size batch_size]  
  [--storage-type mount|remote|S3]  
  [--remote-port port]  
  [--remote-user username]  
  [-s|--shard shard-name]
```

```
[--s3-config-path path]  
[-y|--yes]
```

Adds/removes and enables/disables the archive command for every replication group in the Postgres Pro Shardman cluster to put WAL logs into the initialized backup repository.

add

Adds and enables the `archive` command for every replication group in the Postgres Pro Shardman cluster.

rm

Disables the `archive` command in every replication group in the Postgres Pro Shardman cluster. No additional options are required.

-B *path*

--backup-path *path*

Required when adding `archive_command`. Specifies the path to the backup catalog where the Postgres Pro Shardman cluster backups should be stored.

--batch-size *batch_size*

To speed up the archiving, specify the `--batch-size` option to copy the WAL segments in batches of a specified size. If the `--batch-size` option is used, it is also possible to specify the `-j` option to copy a batch of the WAL segments on multiple threads.

--jobs *count*

-j *count*

The number of parallel threads that `pg_probackup` uses when creating a backup. Default: 1.

--compress

Enables backup compression. If this flag is not specified, compression will be disabled. If the flag is specified, the default `zstd` algorithm is used with the compression level set to 1, while other compression options are ignored even if they are specified.

--compress-algorithm *algorithm*

Defines the compression algorithm: `zlib`, `lz4`, `zstd`, `pglz`, or `none`. Once defined, it checks if the values are valid within the scale of the defined algorithm.

The supported compression algorithms depend on the version of Postgres Pro Shardman that includes the `pg_probackup` used, as explained in [Compression Options](#).

Default: `none`.

--compress-level *level*

Defines the compression level — 0-9 for `zlib`, 1 for `pglz`, 0-22 for `zstd`, and 0-12 for `lz4`.

Default: 1.

--storage-type `mount|remote|S3`

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively.

Default: `remote`.

`--remote-port port`

Specifies the remote ssh port for replication group instances.

Default: 22.

`--remote-user username`

Specifies the remote ssh user for replication group instances.

Default: postgres.

`-s|--shard shard-name`

Specifies the name of the shard where the `archive` command must be added, enabled or disabled. If not specified, the `archive` command is enabled or disabled for every shard.

`--s3-config-path path`

Specifies the path to the S3 configuration file.

Default: <shardman-data-dir>/s3.config.

`-y`

`--yes`

Confirm the restart instead of asking approval from the standard input. Only applies for the `add` command.

backup

Syntax:

```
shardmanctl [common_options] probackup backup -B|--backup-path path
-E|--etcd-path path
-b|--backup-mode MODE
[-j|--jobs count]
[--compress]
[--compress-algorithm algorithm]
[--compress-level level]
[--batch-size batch_size]
[--storage-type mount|remote|S3]
[--remote-port port]
[--remote-user username]
[--ssh-key path]
[-t|--timeout seconds]
[-m|--maxtasks number_of_tasks]
[--log-directory path]
[--s3-config-path path]
[--no-validate]
[--skip-block-validation]
[--log-to-console]
[--retention-redundancy]
[--retention-window]
[--wal-depth]
[--delete-wal]
[--delete-expired]
[--merge-expired]
[-y | --yes]
[--lock-lifetime]
[-e |--external-dirs dir_name]
```

Creates a backup of the Postgres Pro Shardman cluster.

`-B path`

`--backup-path path`

Required. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

`-E path`

`--etcd-path path`

Required. Specifies the path to the catalog where the etcd dumps should be stored.

`-b MODE`

`--backup-mode MODE`

Required. Defines the backup mode: FULL, PAGE, DELTA, PTRACK.

`--batch-size batch_size`

To speed up the archiving, specify the `--batch-size` option to copy the WAL segments in batches of a specified size. If the `--batch-size` option is used, it is also possible to specify the `-j` option to copy a batch of the WAL segments on multiple threads.

`--jobs count`

`-j count`

The number of parallel threads that `pg_probackup` uses when creating a backup. Default: 1.

`--compress`

Enables backup compression. If this flag is not specified, compression will be disabled. If the flag is specified, the default `zstd` algorithm is used with the compression level set to 1, while other compression options are ignored even if they are specified.

`--compress-algorithm algorithm`

Defines the compression algorithm: `zlib`, `lz4`, `zstd`, `pglz`, or `none`.

The supported compression algorithms depend on the version of Postgres Pro Shardman that includes the `pg_probackup` used, as explained in [Compression Options](#).

Default: `none`.

`--compress-level level`

Defines the compression level — 0-9 for `zlib`, 1 for `pglz`, 0-22 for `zstd`, and 0-12 for `lz4`.

Default: 1.

`--remote-port port`

Specifies the remote ssh port for replication group instances.

Default: 22.

`--remote-user username`

Specifies the remote ssh user for replication group instances.

Default: `postgres`.

`--ssh-key path`

Specifies the ssh private key for execution of remote ssh commands.

Default: `$HOME/.ssh/id_rsa`.

`-t seconds`

`--timeout seconds`

Exit with error after waiting until the cluster is ready for the specified number of seconds.

`-m number_of_tasks`

`--maxtasks number_of_tasks`

Specifies the maximum number of concurrent tasks (`pg_probackup` commands) to run.

Default: number of logical CPUs of the system.

`--no-validate`

Skip automatic validation after the backup is taken. You can use this flag if you validate backups regularly and would like to save time when running backup operations.

Default: `false`.

`--skip-block-validation`

Disables block-level checksum verification to speed up the backup process.

Default: `false`.

`--storage-type mount|remote|S3`

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively.

Default: `remote`.

`--log-to-console`

Enables output of the `pg_probackup` logs to the console.

Default: `false`.

`--log-directory path`

Specifies the directory for `pg_probackup` logs. Required if `--storage-type` is set to `S3` unless the `SDM_LOG_DIRECTORY` environment variable is set.

Default: `<backup-directory>/backup/log`.

`--s3-config-path path`

Specifies the path to the S3 configuration file.

Default: `<shardman-data-dir>/s3.config`.

`--retention-redundancy=redundancy`

Specifies the number of full backup copies to keep in the data directory. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--retention-window=window`

Number of days of recoverability. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--wal-depth=wal_depth`

Number of latest valid backups on every timeline that must retain the ability to perform PITR. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--delete-wal`

Deletes WAL files that are no longer required to restore the cluster from any of the existing backups.

Default: `false`.

`--delete-expired`

Deletes backups that do not conform to the retention policy.

Default: `false`.

`--merge-expired`

Merges the oldest incremental backup that satisfies the requirements of retention policy with its parent backups that have already expired.

Default: `false`.

`-y`

`--yes`

Confirm the restart instead of asking approval from the standard input.

`--lock-lifetime`

Allows setting the maximum time that `probackup` can hold the lock, in seconds.

Default: 1800.

`-e dir_name`

`--external-dirs`

Specifies directories outside `PGDATA` to be included in the backup.

checkdb

Syntax:

```
shardmanctl [common_options] probackup checkdb
[--amcheck [--skip-block-validation] [--heapallindexed]] [--shard shard]
[-m|--maxtasks number_of_tasks]
```

Verifies the Postgres Pro Shardman cluster correctness by detecting physical and logical corruption.

`--amcheck`

Performs logical verification of indexes if no corruption was found while checking data files. You must have the [amcheck](#) extension or the [amcheck_next](#) extension installed in the database to check its indexes. For databases without `amcheck`, index verification will be skipped. The `amcheck` extension is included with the Postgres Pro Shardman package.

`--heapallindexed`

Checks that all heap tuples that should be indexed are actually indexed. You can use this flag only together with the `--amcheck` flag. This option is effective depending on the version of `amcheck/am`

check_next installed. The amcheck extension included in the Postgres Pro Shardman package supports this verification.

`--skip-block-validation`

Skip validation of data files. You can use this flag only together with the `--amcheck` flag, so that only logical verification of indexes is performed.

`--shard shard`

Perform the verification only on the specified shard. By default, the verification is performed on all shards.

`-m number_of_tasks`

`--maxtasks number_of_tasks`

Specifies the maximum number of concurrent tasks (pg_probackup commands) to run.

Default: number of logical CPUs of the system.

delete

Syntax:

```
shardmanctl [common_options] probackup delete -B|--backup-path path
[-i|--backup-id backup_id]
[-j|--jobs count]
[-m|--maxtasks number_of_tasks]
[--storage-type mount|remote|S3]
[--s3-config-path path]
[--delete-wal]
[-y|--yes]
[--retention-redundancy]
[--retention-window]
[--wal-depth]
[--delete-expired]
[--merge-expired]
```

Deletes a backup of the Postgres Pro Shardman cluster with specified `backup_id` or launches the retention purge of backups and archived WAL that do not satisfy the current retention policies.

Note that `backup_id` cannot be used with `merge-expired` or `delete-expired`.

`-B path`

`--backup-path path`

Required. Specifies the path to the backup catalog (or key in the bucket of the S3-compatible storage) where Postgres Pro Shardman cluster backups should be stored.

`-i backup_id`

`--backup-id backup_id`

Specifies the unique identifier of the backup.

`--jobs count`

`-j count`

The number of parallel threads that pg_probackup uses when creating a backup. Default: 1.

`-m number_of_tasks`

`--maxtasks number_of_tasks`

Specifies the maximum number of concurrent tasks (pg_probackup commands) to run.

Default: number of logical CPUs of the system.

`--storage-type mount|remote|S3`

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively.

Default: `remote`.

To delete the backup that was created with a `--storage-type` option with a `S3` value, set a `--storage-type` option to a `S3` value in the `delete` command.

`--s3-config-path path`

Specifies the path to the S3 configuration file.

Default: `<shardman-data-dir>/s3.config`.

`--delete-wal`

Deletes WAL files that are no longer required to restore the cluster from any of the existing backups.

Default: `false`.

`-y`

`--yes`

Approve operation.

Default: `false`.

`--retention-redundancy=redundancy`

Specifies the number of full backup copies to keep in the data directory. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--retention-window=window`

Number of days of recoverability. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--wal-depth=wal_depth`

Number of latest valid backups on every timeline that must retain the ability to perform PITR. Must be a non-negative integer. The zero value disables this setting.

Default: current value of the `pg_probackup.conf` file, 0 if not specified.

`--delete-expired`

Deletes backups that do not conform to the retention policy.

Default: `false`.

`--merge-expired`

Merges the oldest incremental backup that satisfies the requirements of retention policy with its parent backups that have already expired.

Default: `false`.

merge

Syntax:

```
shardmanctl [common_options] probackup merge -B|--backup-path path
-i|--backup-id backup_id
[-j|--jobs count]
[-m|--maxtasks number_of_tasks]
[--no-validate]
[--no-sync]
[-y|--yes]
```

Merges the backups that belong to a common incremental backup chain. The full backup merges the backups with their first incremental backup. The incremental backup merges the backups with their parent full backup, along with all the incremental backups between them. Once the merge is complete, the full backup covers all the merged data, and the incremental backups are removed as redundant.

-B *path*
--backup-path *path*

Required. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

-i *backup_id*
--backup-id *backup_id*

Required. Specifies the unique identifier of the backup.

--jobs *count*
-j *count*

The number of parallel threads that pg_probackup uses when creating a backup. Default: 1.

-m *number_of_tasks*
--maxtasks *number_of_tasks*

Specifies the maximum number of concurrent tasks (pg_probackup commands) to run.

Default: number of logical CPUs of the system.

--no-sync

Do not sync merged files to disk. You can use this flag to speed up the merge process. Using this flag can result in data corruption in case of operating system or hardware crash.

Default: false.

--no-validate

Skip automatic validation before and after merge.

Default: false.

-y
--yes

Approve the operation.

Default: false.

restore

Syntax:

```
shardmanctl [common_options] probackup restore
-B|--backup-path path
-i|--backup-id id
-j|--jobs count
```

```
[--recovery-target-time timestamp]  
[-I|--recovery-mode incremental_mode]  
[-t|--timeout seconds]  
[-m|--maxtasks number_of_tasks]  
[--metadata-only] [--schema-only] [--shard shard]  
[--no-validate]  
[--skip-block-validation]  
[--s3-config-path path]  
[--storage-type mount|remote|S3]  
[--wal-limit number_of_wal_segments]  
[--log-directory path]  
[--data-validate]  
[-e|--external-dirs dir_name]  
[-T|--tablespace-mapping old_dir=new_dir]  
[-e|--external-mapping old_dir=new_dir]  
[-skip-external-dirs]  
[--reinit-before-restore]
```

Restores a Postgres Pro Shardman cluster from the selected backup.

`-B path`
`--backup-path path`

Required. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

`-i id`
`--backup-id id`

Required. Specifies backup ID for restore.

`--jobs count`
`-j count`

The number of parallel threads that `pg_probackup` uses when restoring from a backup. Default: 1.

`--recovery-target-time timestamp`

Point-in-Time Recovery (PITR) option. Specifies the timestamp for restore. Example: '2024-01-25 15:30:36' in UTC.

`-I incremental_mode`
`--recovery-mode incremental_mode`

Specifies the incremental restore mode to be used. Possible values are:

- `checksum` — replace only pages with mismatched checksum and LSN.
- `lsn` — replace only pages with LSN greater than point of divergence.
- `none` — regular restore, default.

`-t seconds`
`--timeout seconds`

Exit with error after waiting until the cluster is ready or the recovery is complete for the specified number of seconds.

`--metadata-only`

Perform metadata-only restore. By default, full restore is performed.

`--schema-only`

Perform schema-only restore. By default, full restore is performed.

`--shard shard`

Perform restoring only on the specified shard. By default, restoring is performed on all shards.

`--no-validate`

Skip backup validation. You can use this flag if you validate backups regularly and would like to save time when running restore operations.

Default: `false`.

`--skip-block-validation`

Disable block-level checksum verification to speed up validation. During automatic validation before the restore only file-level checksums will be verified.

Default: `false`.

`--s3-config-path path`

Specifies the path to the S3 configuration file.

Default: `<shardman-data-dir>/s3.config`.

`--storage-type mount|remote|S3`

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively. When creating backup with a `--storage-type` option with a `S3` value, set `--storage-type` option to a `S3` value in the `restore` command.

Default: `remote`.

`--wal-limit number_of_wal_segments`

Specifies the number of WAL segments in which the closest synchronization points will be searched in the case of PITR.

Default: `0` — no limit.

`--log-directory path`

Specifies the directory for `pg_probackup` logs. Required if `--storage-type` is set to `S3` unless the `SDM_LOG_DIRECTORY` environment variable is set.

Default: `<backup-directory>/backup/log`.

`--data-validate`

If enabled, verifies data with `probackup validate` before restoring.

Default: `false`.

`-e dir_name`

`--external-dirs`

Specifies directories outside `PGDATA` to be included in the backup.

`-T`

`--tablespace-mapping old_dir{rgid}=new_dir{rgid}`

Specifies mapping for tablespaces. When being restored, the content of `old_dir` is moved to `new_dir`, if the tablespaces still exist and their location corresponds. Like with [CREATE TABLESPACE](#), you can use `{rgid}` value for a replication group ID.

`-e`
`--external-mapping old_dir=new_dir`
Specifies mapping for external directories from the `--external-dirs` parameter. When being re-stored, the content of *old_dir* is moved to *new_dir*, if the directories still exist.

`--skip-external-dirs`
If specified, no external directories are restored.

`--reinit-before-restore`
Reinitialize the cluster with the same topology before restore.
Default: false.

show**Syntax:**

```
shardmanctl [common_options] probackup show
  -B|--backup-path path
  [-f|--format table|json]
  [--archive ]
  [-i|--backup-id backup-id]
  [--instance instance]
  [--storage-type mount|remote|S3]
  [--s3-config-path path]
```

Shows the list of backups of the Postgres Pro Shardman cluster.

`-B path`
`--backup-path path`
Required. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

`-f table|json`
`--format table|json`
Specifies the output format.
Default: `table`.

`--archive`
Shows the WAL archive information.

`-i backup-id`
`--backup-id backup-id`
Shows information about the specific backups.

`--instance instance`
Shows information about the specific instance.

`--s3-config-path path`
Specifies the path to the S3 configuration file.
Default: `<shardman-data-dir>/s3.config`.

`--storage-type mount|remote|S3`
Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object

storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively. To show a backup that was created with the `S3` value of `--storage-type`, set `--storage-type` to `S3` in the `show` command.

Default: `remote`.

show-config

Syntax:

```
shardmanctl [common_options] probackup show-config
  -B backup_path
  [--format=text|json]
  [--no-scale-units]
  -s|--shard shard_name
  [--s3-config-path path]
  [--storage-type mount|remote|S3]
```

Displays all the current `pg_probackup` configuration settings, including those that are specified in the `pg_probackup.conf` configuration file located in the `backup_dir/backups/shard_name` directory and those that were provided on a command line.

`-B string`
`--backup-path=string`

Required. Specifies the absolute path to the backup catalog.

`--format text|json`

Specifies the output format.

Default: `text`.

`--no-scale-units`

Output the configuration parameter values for the time and the amount of memory in the default units.

Default: `false`.

`-s string`
`--shard=string`

A name of the shard to execute the `show-config` command for.

`--s3-config-path path`

Specifies the path where the S3 configuration file will be created on all Postgres Pro Shardman nodes.

Default: `<shardman-data-dir>/s3.config`.

`--storage-type mount|remote|S3`

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively.

Default: `remote`.

validate

Syntax:

```
shardmanctl [common_options] probackup validate
  -B|--backup-path path
```

```
-i|--backup-id id
[-t|--timeout seconds]
[-m|--maxtasks number_of_tasks]
[--log-to-console]
[--storage-type mount|remote|S3]
[--s3-config-path path]
[--log-directory path]
[--remote-port port]
[--remote-user username]
```

Checks the selected Postgres Pro Shardman cluster backup for integrity.

```
-B path
--backup-path path
```

Required. Specifies the path to the backup catalog where Postgres Pro Shardman cluster backups should be stored.

```
-i id
--backup-id id
```

Required. Specifies backup ID for validation.

```
--log-to-console
```

Enables output of `pg_probackup` logs to the console.

Default: `false`.

```
-t seconds
--timeout seconds
```

Exit with error after waiting until the cluster is ready for the specified number of seconds.

```
-m number_of_tasks
--maxtasks number_of_tasks
```

Specifies the maximum number of concurrent tasks (`pg_probackup` commands) to run.

Default: number of logical CPUs of the system.

```
--s3-config-path path
```

Specifies the path to the S3 configuration file.

Default: `<shardman-data-dir>/s3.config`.

```
--storage-type mount|remote|S3
```

Type of the backup storage. If the value is `remote`, SSH is used to copy data files to the remote backup directory. But this behavior is different if a directory mounted to all nodes or an S3-compatible object storage is used to store backups. To specify these kinds of storage, the value of the `--storage-type` option is set to `mount` or `S3`, respectively. To validate a backup that was created with the `S3` value of `--storage-type`, set `--storage-type` to `S3` in the `validate` command.

Default: `remote`.

```
--log-directory path
```

Specifies the directory for `pg_probackup` logs. Required if `--storage-type` is set to `S3` unless the `SDM_LOG_DIRECTORY` environment variable is set.

Default: `<backup-directory>/backup/log`.

`--remote-port port`

Specifies the remote ssh port for replication group instances.

Default: 22.

`--remote-user username`

Specifies the remote ssh user for replication group instances.

Default: postgres.

`--ssh-key path`

Specifies the ssh private key for execution of remote ssh commands.

Default: `$HOME/.ssh/id_rsa`.

set-config

Syntax:

```
shardmanctl [common_options] probackup set-config
  [--archive-timeout int]
  [-B | --backup-path string]
  [-m | --maxtasks int]
  [--remote-port int]
  [--remote-user string]
  [--retention-redundancy int]
  [--retention-window int]
  [--wal-depth int]
  [--s3-config-path string]
  [-s | --shard string]
  [--storage-type string]
  [-e | --external-dirs dir_name]
```

Adds the specified settings to the `pg_probackup.conf` or modifies those previously added.

`--archive-timeout int`

Sets a timeout for the WAL segment archiving and streaming, in seconds.

Default: `pg_probackup` waits for 300 seconds.

`-B string`

`--backup-path=string`

Specifies the absolute path to the backup catalog.

`-m int`

`--maxtasks=int`

Specifies the maximum number of concurrent tasks (`pg_probackup` commands) to run.

Default: number of logical CPUs of the system.

`--remote-port int`

An SSH remote backup port.

Default: 22.

`--remote-user string`

An SSH remote backup user.

`--retention-redundancy int`

Specifies the number of the full backup copies to store in the data directory. It must be set to a non-negative integer. The zero value disables this setting.

Default: 0.

`--retention-window int`

A number of days of recoverability. It must be set to a non-negative integer. The zero value disables this setting.

Default: 0.

`--wal-depth int`

A number of the latest valid backups on every timeline that must retain the ability to perform PITR. Must be set to a non-negative integer. The zero value disables this setting.

`--s3-config-path string`

A path to the S3 configuration file.

Default: `/var/lib/pgpro/sdm-17/data/s3.config`

`-s string`

`--shard=string`

A name of the shard to make the `set-config` command for. If not specified, the command is run for all the shards.

Default: current value of the `pg_probackup.conf` file.

`--storage-type string`

A backup storage type, the possible values are `remote`, `mount`, `S3`.

Default: `remote`.

`-e dir_name`

`--external-dirs`

Specifies directories outside `PGDATA` to be included in the backup.

Backing up a Postgres Pro Shardman Cluster Using `probackup` Command

On the backup host:

```
groupadd postgres
useradd -m -N -g postgres -r -d /var/lib/postgresql -s /bin/bash
```

Then add SSH keys to provide passwordless SSH connection between the backup host and Postgres Pro Shardman cluster hosts. Then on the backup host:

```
apt-get install pg-probackup shardman-utils
mkdir -p directory
chown postgres:postgres directory -R
shardmanctl [common_options] probackup init --backup-path=directory --etcd-
path=directory/etcd --remote-user=postgres --remote-port=22
shardmanctl [common_options] probackup archive-command --backup-path=directory --
remote-user=postgres --remote-port=22
```

If all the requirements are met, then run the `backup` subcommand for the cluster backup:

```
shardmanctl [common_options] probackup backup --backup-path=directory --etcd-  
path=directory --backup-mode=MODE
```

You must pass the directories through the `--backup-path` and `--etcd-path` options and backup mode through `--backup-mode`. Full and delta backups are available with `FULL`, `DELTA`, `PTRACK` and `PAGE` values. Also it is possible to specify backup compression options through `--compress`, `--compress-algorithm` and `--compress-level` flags, as well as specify `--remote-port` and `--remote-user` flags. You can limit the number of running concurrent tasks when doing backup by passing the limit through the `--maxtasks` flag.

By default, copying data via SSH is used to create a backup. To copy data to a mounted partition instead, use the `--storage-type` option with the `mount` value. This value will be automatically used in the restore process.

You can also copy data to an S3-compatible object storage. To do this, use the `--storage-type` option with the `S3` value. When this value is used, it is required to specify the directory for `pg_probackup` logs. You can do it either by specifying `--log-directory` for each command or set the environment variable `SDM_LOG_DIRECTORY`, for example:

```
export SDM_LOG_DIRECTORY=/backup/logs
```

If you are going to perform backup/restore only for an S3-compatible object storage, you can also set an environment variable instead of specifying `--storage-type` in each `probackup` command:

```
export SDM_STORAGE_TYPE=S3
```

Restoring a Postgres Pro Shardman Cluster using probackup command

`shardmanctl` in `probackup` mode can perform either full restore, metadata-only or schema-only restore of a Postgres Pro Shardman cluster from a backup created by the `probackup backup` command.

To perform full or partial restore, firstly you must select needed backup to restore from. To show list of available backups run the following command:

```
shardmanctl [common_options] probackup show --backup-path=path --format=format [--  
archive ] [-i|--backup-id backup-id] [--instance instance]
```

The output should be a list of backups with their IDs in a table or JSON format. Then pick the needed backup ID and run the `probackup restore` command.

```
shardmanctl [common_options] probackup restore --backup-path=path --backup-id=id
```

Pass the path to the repo through the `--backup-path` option and backup ID through `--backup-id` flag.

If you encounter issues with an `etcd` instance, it makes sense to perform metadata-only restore. To do this, you can run the following command:

```
shardmanctl [common_options] probackup restore --backup-path=path --backup-id=id --  
metadata-only
```

If you need to restore only schema information, like: tables, roles and etc. you should specify `--schema-only` option.

For both kinds of restore, you can specify `--timeout` for the tool to exit with error after waiting until the cluster is ready or the recovery is complete for the specified number of seconds.

You can specify `--shard` parameter for restoring only on the single shard.

Also you can specify `--recovery-target-time` option for Point-in-Time Recovery. In this case Postgres Pro Shardman finds the closest syncpoint to specified timestamp and suggests restoring on the found LSN. You can also specify `--wal-limit` to limit the number of WAL segments to be processed.

Important

Before restoring a Postgres Pro Shardman cluster, make sure that the cluster is up by executing the `shardmanctl status` command. If the output shows errors, performing the restore can result in the cluster becoming unavailable. First, fix the errors by reinitializing the cluster and restoring the etcd metadata. Then you can proceed to restoring the cluster from backup.

Restoring a Postgres Pro Shardman Cluster

`shardmanctl` can perform either full restore, metadata-only or schema-only restore of a Postgres Pro Shardman cluster from a backup created by the `backup` command.

To perform full restore, you can run the following command:

```
shardmanctl [common_options] recover --info file
```

Pass the file to load information about the backup from through the `--info` option. In most cases, set this option to point to the `backup_info` file in the backup directory or to its modified copy.

If you encounter issues with an etcd instance, it makes sense to perform metadata-only restore. To do this, you can run the following command:

```
shardmanctl [common_options] recover --dumpfile file --metadata-only
```

You must pass the file to load the etcd metadata dump from through the `--dumpfile` option.

If you need to restore only schema information, like: tables, roles and etc. you should specify `--schema-only` option.

For all kinds of restore, you can specify `--timeout` for the tool to exit with error after waiting until the cluster is ready or the recovery is complete for the specified number of seconds.

You can specify `--shard` parameter for restoring only on the single shard.

Before running the `recover` command, specify `DataRestoreCommand` and `RestoreCommand` in the `back-up_info` file. `DataRestoreCommand` fetches the base backup and restores it to the shard data directory. `RestoreCommand` fetches the WAL file and saves it to the shard `pg_wal` directory. These commands can use the following substitutions:

`%p`

Destination path on the server.

`%s`

SystemId of the restored database (the same in the backup and in restored cluster).

`%f`

Name of the WAL file to restore.

Shard `keeper` thread runs both commands on each node in the cluster. Therefore:

- Make the backup accessible to these nodes (for example, by storing it in a shared filesystem or by using a remote copy protocol, such as SFTP).
- Commands to fetch the backup are executed as the operating system user under which shards work (usually `postgres`), so set the permissions for the backup files appropriately.

These examples show how to specify `RestoreCommand` and `DataRestoreCommand`:

- If a backup is available through a passwordless SCP, you can use:


```
"DataRestoreCommand": "scp -r user@host:/var/backup/shardman/%s/backup/* %p",  
"RestoreCommand": "scp user@host:/var/backup/shardman/%s/wal/%f %p"
```

- If a backup is stored on NFS and available through `/var/backup/shardman path`, you can use:

```
"DataRestoreCommand": "cp -r /var/backup/shardman/%s/backup/* %p",  
"RestoreCommand": "cp /var/backup/shardman/%s/wal/%f %p"
```

Configuration-related Commands

config get

Syntax:

```
shardmanctl [common_options] config get [-c | --choose-revision] [-r | --revision ] [-f  
| --file]
```

Outputs the current full cluster specification or a configuration of the specified revision.

```
-c  
--choose-revision
```

Enables an interactive mode of choosing a configuration of the specified revision.

```
-r  
--revision
```

ID of a configuration revision.

```
-f file_name  
--file=file_name
```

Name of a file for writing the configuration. If not specified, the value is stdout.

config revisions rm

Syntax:

```
shardmanctl [common_options] config revisions rm [-r | --revision ] [-y | --yes]
```

Deletes a specified configuration revision from history.

```
-r  
--revision
```

ID of a configuration revision. If not specified, enables an interactive mode of choosing a configuration of the specified revision. This is a timestamp of an operation that resulted in Postgres Pro Shardman configuration change.

```
-y  
--yes
```

Perform automatic confirmation.

config update

Syntax:

```
shardmanctl [common_options] config update [[-f|--file shard_spec_file /  
shardman_spec_file]|spec_text [-p|--patch][-w|--wait]] [--force] [-y | --yes]
```

Updates the shard or full Postgres Pro Shardman configuration. Supports the following placeholders for postgres parameters: `{{cluster}}`, `{{shard}}`, `{{dataDir}}`, `{{keeperDir}}`, `{{keeperID}}`, `{{host}}`. See example [here](#).

```
-f shard_spec_file | shardman_spec_file
--specfile=shard_spec_file | shardman_spec_file
```

Specifies the file with the shard or full Postgres Pro Shardman configuration. The configuration file type is determined automatically. The value of `-` means the standard input. By default, the configuration is passed in *spec_text*.

```
-w
--wait
```

Sets `shardmanctl` to wait for configuration changes to take effect. If a new configuration cannot be loaded by all replication groups, `shardmanctl` will wait forever.

```
-p
--patch
```

Merge the new configuration into the existing one. By default, the new configuration replaces the existing one.

```
--force
```

Perform forced update if a cluster operation is in progress.

```
-y
--yes
```

Confirm the restart necessary for the parameters to take effect. If this option is not specified, and the parameters update requires a restart, the manual confirmation will be requested. If not confirmed, the cluster will continue to work, yet the new parameter values will only take effect after the restart.

config rollback

Syntax:

```
shardmanctl [common_options] config rollback [-r | --revision] [-w|--wait time_duration] [--force] [-y|--yes][--repair]
```

Makes a rollback of Postgres Pro Shardman to one of the previous states. When rolling back to the config revision that has `max_connections`, `max_prepared_transactions`, or `max_worker_processes` parameters, the replicas are reinitialized.

When rolling back to one of the previous states, none of the `max_connections`, `max_prepared_transactions`, `max_locks_per_transaction`, `max_wal_senders`, `max_worker_processes` parameters is rolled back if its current value is higher than the one it is rolled back to.

```
-r
--revision
```

ID of a revision the rollback must be made to. It is a timestamp of an operation that resulted in Postgres Pro Shardman configuration change.

If not specified, a user is presented with a list of revisions that he can choose from.

```
-w
--wait
```

Sets `shardmanctl` to wait for configuration changes to take effect. If a new configuration cannot be loaded by all replication groups, `shardmanctl` will wait forever.

Default: 1h.

```
-f
--force
```

Perform forced setting of a parameter if a cluster operation is in progress.

`-y`
`--yes`

Perform automatic confirmation.

`--repair`

Enables the mode of restoring the cluster that failed due to the invalid parameter values.

Default: disabled.

config revisions

Syntax:

```
shardmanctl [common_options] config revisions [-f|--format text|json]
```

Outputs the revision history of the Postgres Pro Shardman cluster configuration. It has the following information for each revision:

- `revision_id` — timestamp of the command that resulted in the Postgres Pro Shardman cluster configuration change
- `host` — name of the host from which this command was executed
- `user` — user who executed this command
- `command` — the command itself

`-f=text|json`
`--format=text|json`

Specifies the output format.

Default: `text`.

config revisions set

Syntax:

```
shardmanctl [common_options] config revisions set [--keep-config-revisions]
```

Allows setting the length of the configuration revision history. This length cannot be lower than 5, in which case it is automatically set to 5. For Postgres Pro Shardman clusters where the configuration revision history was not collected yet, the length is automatically set to 20.

`--keep-config-revisions`

A limit on the number of revisions for one Postgres Pro Shardman configuration. If the limit is lower than the current history length, the older versions out of this limit will be deleted. Also, if the number of operations resulting in configuration changes exceeds the limit, the oldest revision is deleted.

Default: 20.

config update ip

Syntax:

```
shardmanctl [common_options] config update ip [-u|ip_1=ip_2,hostname_1=hostname_2] [-y|--yes] [-repair]
```

Updates the specified node IPs in the cluster.

`-u`
`ip_1=ip_2,hostname_1=hostname_2`

Specifies the node IPs to be updated.

`-y`
`--yes`

Perform automatic confirmation.

`--repair`

Enables the mode of restoring the cluster that failed due to the invalid parameter values.

Default: disabled.

config update credentials

Syntax:

```
shardmanctl [common_options] config update credentials [-u | --user] [-p | --password]
[-k | --ssl-key] [-c | --ssl-cert] [-w|--wait time_duration] [--force] [-y | --yes]
```

Updates password or certificate/key of a user to connect to a Postgres Pro Shardman cluster. It only updates the authentication type that was specified by the user (`scram-sha-256`, `ssl`) and not the type itself.

`-u`
`--user`

User that requires an update of the authentication parameters.

`-p`
`--password`

New password.

`-k`
`--ssl-key`

New SSL key.

`-c`
`--ssl-cert`

New SSL certificate.

`-w`
`--wait`

Sets `shardmanctl` to wait for configuration changes to take effect. If a new configuration cannot be loaded by all replication groups, `shardmanctl` will wait forever.

`--force`

Perform forced update if a cluster operation is in progress.

`-y`
`--yes`

Confirm the operation instead of asking approval from the standard input.

Cluster-related Commands

cluster repfactor set

Syntax:

```
shardmanctl [common_options] cluster repfactor set --value new_repfactor
```

Sets the replication factor for the manual-topology mode.

`--value=new_repfactor`

New replication factor value

cluster start

Syntax:

`shardmanctl [common_options] cluster start`

Starts all PostgreSQL server instances.

cluster stop

Syntax:

`shardmanctl [common_options] cluster stop [-y|--yes] [-m|--message]`

Stops all PostgreSQL server instances.

`-y`

`--yes`

Confirm the operation instead of asking approval from the standard input.

`-m`

`--message`

Explains why the cluster was stopped.

cluster topology

Syntax:

`shardmanctl [common_options] cluster topology -f|--format table|json|text`

Displays the cluster topology.

`-f table|json|text`

`--format=table|json|text`

Output format. For more details, see [the section called “Displaying the Cluster Topology”](#).

init

Syntax:

`shardmanctl [common_options] init [-y|--yes] [-f|--spec-file spec_file_name]|spec_text`

Registers a new Postgres Pro Shardman cluster in the etcd store or reinitializes the existing cluster defining a new cluster configuration and removing all data and nodes.

`-f spec_file_name`

`--specfile=spec_file_name`

Specifies the file with the cluster specification string. The value of `-` means the standard input. By default, the string is passed in `spec_text`. For usage details, see [the section called “Registering a Postgres Pro Shardman Cluster”](#).

`-y`

`--yes`

Confirm the operation instead of asking approval from the standard input.

intcheck

Syntax:

```
shardmanctl [common_options] intcheck [-s|--system] [-u|--user] [-c|--catalog] [-o|--output] [-n|--node node]
```

Runs `pg_integrity_check` on all nodes of a Postgres Pro Shardman cluster or on a selected one node.

`-s`
`--system`
Validate checksums for read-only files. Checksums for read-only files control both file contents and file attributes.

`-u`
`--user`
Validate checksums for additional files. Checksums for additional files control both file contents and file attributes.

`-c`
`--catalog`
Validate checksums for system catalog tables. For the `-c` option to work correctly, the database server must be started and accept connections.

`-o`
`--output`
Recalculate checksums and write them into a file

`-n node_names`
`--node=node_names`
Only execute the `pg_integrity_check` command on the selected node

load

Syntax:

```
shardmanctl [common_options] load [ -b | --batch-size lines_limit] [ --destination-fields fields_list ]  
[ --distributed-keys key_type_list] [ -D | --delimiter character ]  
[ --null-marker string] [ -e | --escape character] [ -f | --file input_file ]  
[ -F | --format text/csv ] [ -j | --jobs task_total] [ -q | --quote character ]  
[ --reject-file filename] [ --schema filename] [ --source file/postgres ]  
[ --source-connstr connect_string] [ --source-fields fields_list] [ --source-table source_table ]  
[ -t | --table destination_table ]
```

Loads data to a Postgres Pro Shardman cluster.

`-b lines_limit`
`--batch-size=lines_limit`
Number of rows per batch to write to the Postgres Pro Shardman cluster.

Default: 1000.

`--destination-fields=fields_list`
Comma-separated list of target table fields. If the value is not set, then all fields of the table are used in the order they are declared.

`--distributed-keys=key_type_list`

Comma-separated list of pairs. Each pair consists of a field number (starting with zero) and a type, which are separated by a colon. The following types are supported: `bool`, `char`, `float4`, `float8`, `int2`, `int4`, `int8`, `name`, `text`, `varchar` and `uuid`.

`-D character`

`--delimiter=character`

Specifies the character that separates columns within each row (line) of the file. This must be a single one-byte character.

Default: `tab` for text format, `comma` for CSV format

`--null-marker=string`

Specifies the string that represents a null value.

Default: `\N` for text format, unquoted empty string for CSV format.

`-e character`

`--escape=character`

Specifies the character that should appear before a data character that matches the `QUOTE` value. The default is the same as the `QUOTE` value (so that the quoting character is doubled if it appears in the data). This must be a single one-byte character. This option is allowed only when using CSV format.

`-f filename`

`--file=filename`

Input data filename (or `-` for stdin)

`-F text/csv`

`--format=text/csv`

Input data format. Possible values are `text` and `csv`.

Default: `text`.

`-j number`

`--jobs=number`

Number of parallel processes to load data.

Default: number of replication groups.

`-q character`

`--quote=character`

Specifies the quoting character to be used when a data value is quoted. The default is double-quote. This must be a single one-byte character. This option is allowed only when using CSV format.

`--reject-file=filename`

All data batches with errors during upload will be written to this file. If the value is not set, then such batches will be skipped.

`--schema=filename`

The schema that defines the rules for transferring data from PostgreSQL to Postgres Pro Shardman. If this option is set, then all other options are not used.

`--source=file/postgres`

Data source type — `file` or `postgres`.

Default: `file`.

`--source-connstr=string`

Data source database connection string

`--source-fields=fields_list`

Comma-separated list of source table fields. If the value is not set, then all fields of the table are used in the order they are declared.

`--source-table=table`

Source table, view or function (`funcname (param1, ..., paramN)`).

`-t table`

`--table=table`

Destination table.

restart

Syntax:

```
shardmanctl [common_options] restart [-y|--yes] [--no-wait] [-m|--message]
```

Restarts a Postgres Pro Shardman cluster.

`-y`

`--yes`

Confirm the operation instead of asking approval from the standard input.

`--no-wait`

Do not wait for the replicas to start.

`-m`

`--message`

Explains why the cluster was restarted.

set

Syntax:

```
shardmanctl [common_options] set param1=value1 [param2=value2 [...]] [-y|--yes] [-w|--wait time_duration] [-f|--force] [--repair]
```

Sets the values of the specified Postgres Pro Shardman cluster database parameters. Supports the following placeholders for postgres parameters: `{{cluster}}`, `{{shard}}`, `{{dataDir}}`, `{{keeperDir}}`, `{{keeperID}}`, `{{host}}`. See example [here](#).

`-w`

`--wait`

Sets `shardmanctl` to wait for configuration changes to take effect. Value examples: `2h45m`, `1m30s`, `5m`, `10s`.

Default: `1h`.

`-y`

`--yes`

Confirm the restart necessary for the parameters to take effect. If this option is not specified, and the parameters update requires a restart, the manual confirmation will be requested. If not confirmed, the cluster will continue to work, yet the new parameter values will only take effect after the restart.

`-f`
`--force`

Perform forced setting of a parameter if a cluster operation is in progress.

`--repair`

Enables the mode of restoring the cluster that failed due to the invalid parameter values.

Default: disabled.

unset

Syntax:

```
shardmanctl [common_options] unset pgParam1 [pgParam2 [...]] [-y|--yes] [-w|--wait time_duration] [-f|--force] [--repair]
```

Resets the specified PostgreSQL parameters to their default values.

`-w`
`--wait`

Sets the maximum time that shardmanctl can wait for configuration changes to take effect.

Default: 1h.

`-y`
`--yes`

Confirm the restart necessary for the parameters to take effect. If this option is not specified, and the parameters update requires a restart, the manual confirmation will be requested. If not confirmed, the cluster will continue to work, yet the new parameter values will only take effect after the restart.

`-f`
`--force`

Perform forced setting of a parameter if a cluster operation is in progress.

`--repair`

Enables the mode of restoring the cluster that failed due to the invalid parameter values.

Default: disabled.

Registering a Postgres Pro Shardman Cluster

To register a Postgres Pro Shardman cluster in the etcd store, run the following command:

```
shardmanctl [common_options] init [-y|--yes] [-f|--spec-file spec_file_name] | spec_text
```

You must provide the string with the cluster specification. You can do it as follows:

- On the command line — do not specify the `-f` option and pass the string in `spec_text`.
- On the standard input — specify the `-f` option and pass `-` in `spec_file_name`.
- In a file — specify the `-f` option and pass the filename in `spec_file_name`.

Displaying the Cluster Topology

`cluster topology` displays the current cluster topology. The default is the table mode. All cluster nodes will be grouped by the replication groups they belong to. For each node, its status will be displayed.

```
shardmanctl [common_options] cluster topology -f|--format table|json|text
```

Nodes-related Commands

nodes add

Syntax:

```
shardmanctl [common_options] nodes add -n|--nodes node_names [--no-rebalance]
```

Adds nodes to a Postgres Pro Shardman cluster. Once the command is executed, a dump is created.

```
-n node_names  
--nodes=node_names
```

Required.

Specifies the comma-separated list of nodes to be added.

```
--no-rebalance
```

Skip the step of rebalancing partitions of sharded tables. For more details, see [the section called “Adding Nodes to a Postgres Pro Shardman Cluster”](#).

nodes rm

Syntax:

```
shardmanctl [common_options] nodes rm -n|--nodes node_names
```

Removes nodes from a Postgres Pro Shardman cluster. Once the command is executed, a dump is created.

```
-n node_names  
--nodes=node_names
```

Specifies the comma-separated list of nodes to be removed. For usage details, see [the section called “Removing Nodes from a Postgres Pro Shardman cluster”](#).

nodes start

Syntax:

```
shardmanctl [common_options] nodes start -n|--nodes node_names [--no-wait]
```

Starts the nodes.

```
-n node_names  
--nodes=node_names
```

Node names.

```
--no-wait
```

Sets shardmanctl not to wait for the nodes to start.

nodes restart

Syntax:

```
shardmanctl [common_options] nodes restart -n|--nodes node_names [--no-wait]
```

Restarts the nodes.

```
-n node_names
--nodes=node_names
```

Node names.

```
--no-wait
```

Do not wait for the nodes to restart.

nodes stop

Syntax:

```
shardmanctl [common_options] nodes stop -n|--nodes node_names [--no-wait]
```

Stops the nodes.

```
-n node_names
--nodes=node_names
```

Node names.

```
--no-wait
```

Do not wait for the nodes to stop.

Adding Nodes to a Postgres Pro Shardman Cluster

To add nodes to a Postgres Pro Shardman cluster, run the following command:

```
shardmanctl [common_options] nodes add -n|--nodes node_names
```

You must specify the `-n` (`--nodes`) option to pass the comma-separated list of nodes to be added. Nodes can be referred by their hostname or IP address. Hostnames must be correctly resolved on all nodes.

If `nodes add` command fails during execution, use the `cleanup --after-node-operation` command to fix possible cluster configuration issues.

Removing Nodes from a Postgres Pro Shardman cluster

To remove nodes from a Postgres Pro Shardman cluster, run the following command:

```
shardmanctl [common_options] nodes rm -n|--nodes node_names
```

Specify the `-n` (`--nodes`) option to pass the comma-separated list of nodes to be removed. Recreates all partitions of sharded tables.

Note

Do not use the `cleanup` command to fix possible cluster configuration issues after a failure of `nodes rm`. Redo the `nodes rm` command instead.

To remove all nodes in a cluster and not care about the data, just reinitialize the cluster. If a removed replication group contains local (non-sharded and non-global) tables, the data is silently lost after the replication group removal.

Shard-related Commands

shard add

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name add -n|--nodes node_names [--no-wait]
```

Adds a replica to a shard. Once the command is executed, a dump is created.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
-n node_names  
--nodes=node_names
```

Specifies the comma-separated list of replica nodes to be added.

```
--no-wait
```

Do not wait for the shard to start.

shard master set

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name master set -n| node node_names
```

Sets the precedence for a certain primary server for a specified shard.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
master set
```

Primary server with precedence.

```
-n node_names  
--nodes=node_names
```

Specifies the comma-separated list of replica nodes.

shard master reset

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name master reset
```

Resets the parameters of the master with precedence for the shard.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
master reset
```

Resets the parameters of the master with precedence for the shard.

```
-n node_names  
--nodes=node_names
```

Specifies the comma-separated list of replica nodes.

shard add

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name reset [--yes | -y][--new-  
primary | -p]
```

Resets nodes of a replication group if they are in a state of hanging.

`-s shard_name`
`--shard=shard_name`

Shard name.

`-y`
`--yes`

Confirm the operation instead of asking approval from the standard input.

`--new-primary`
`-p`

New primary node host.

shard rm

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name rm -n|--nodes node_names  
[-f|--force]
```

Removes a replica from a shard. Once the command is executed, a dump is created.

`-s shard_name`
`--shard=shard_name`

Shard name

`-n node_names`
`--nodes=node_names`

Specifies the comma-separated list of replica nodes to be removed.

`-f`
`--force`

Perform forced removal of the node, even if it is dead.

shard switch

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name switch [--new-  
primary node_names]
```

Switches the primary node.

`-s shard_name`
`--shard=shard_name`

Shard name.

`--new-primary=node_names`

New primary node host.

shard start

Syntax:

```
shardmanctl [common_options] shard -s |--shard shard_name start [--no-wait] [-n|--  
node node_name]
```

Starts the shard.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
--no-wait
```

Do not wait for the shard to start.

```
-n node_name  
--node=node_name
```

Specifies the node to start.

shard stop

Syntax:

```
shardmanctl [common_options] shard -s |--shard shard_name stop [-n|--  
node node_name]
```

Stops the shard.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
-n node_name  
--node=node_name
```

Specifies the node to stop.

shard replicas reinit

Syntax:

```
shardmanctl [common_options] shard -s|--shard shard_name replicas reinit [-n|--  
node node_names] [-y|--yes] [--no-wait]
```

Resets replicas of a specific shard.

```
-s shard_name  
--shard=shard_name
```

Shard name.

```
-n node_names  
--node=node_names
```

Specifies the node on which to reset replicas. If not specified, checks shard replicas on all nodes.

```
-y  
--yes
```

Confirm the operation instead of asking approval from the standard input.

```
--no-wait
```

Do not wait wait for replicas to become ready.

For more details, see [the section called “Reinitializing Replicas”](#)

Reinitializing Replicas

If replicas are in an incorrect state, you can reset them using the `shardmanctl` command:

```
shardmanctl [common_options] shard --shard=shard_name replicas reinit
```

This command determines the nodes on which replicas of the specified shard are running and sends a request to `shardmand` on these nodes. After receiving this request, `shardmand` clears the postgres data directory and restarts the `keeper` thread that is responsible for managing the replica. After that, the replicas are restarted and begin to receive data from the corresponding primary.

Status-related Commands

status

Syntax:

```
shardmanctl [common_options] status [-f|--format table|json] [--  
filter store|metadata|shardmand|rg|master|dictionary|biha|keeper|restart_required_params]  
[-s|--sort node|rg|status]  
[--  
exclude store|metadata|shardmand|rg|master|dictionary|biha|keeper|restart_required_params]
```

Reports on the health status of Postgres Pro Shardman cluster subsystems. Shows no details or status tables if no issues are found.

```
-f table|json  
--format=table|json
```

Specifies the report format.

Default: `table`.

For more details, see [the section called “Getting the Status of Cluster Subsystems”](#).

```
--filter          store|metadata|shardmand|rg|master|dictionary|biha|keeper|restart_re-  
quired_params
```

Specifies subsystems whose status information should be included in the output.

Default: `all`.

For more details, see [the section called “Getting the Status of Cluster Subsystems”](#).

```
-s node|rg|status  
--sort node|rg|status
```

Sort messages inside one group (table) as specified.

Default: `node`.

For more details, see [the section called “Getting the Status of Cluster Subsystems”](#).

```
--exclude          store|metadata|shardmand|rg|master|dictionary|biha|keeper|restart_re-  
quired_params
```

Specifies subsystems that should not be included in the `status` output.

status transactions

Syntax:

```
shardmanctl [common_options] status transactions [-r|--  
repgroup replication_group_name]
```

Shows distributed transactions that Postgres Pro Shardman built-in monitoring tools failed to resolve.

```
-r replication_group_name
--repgroup=replication_group_name
```

Specifies the replication group for which to output transactions.

Default: all replication groups.

For more details, see [the section called “Outputting the List of Unresolved Distributed Transactions”](#).

Getting the Status of Cluster Subsystems

To get a report on the health status of Postgres Pro Shardman cluster in a table format for `metadata` and `store` subsystems sorted by replication group, run the following command:

```
shardmanctl [common_options] status --filter=metadata,store --sort=rg
```

To get the report in JSON format, use `-f|--format=json` option (omitted above since `table` format is used by default). Each detected issue is reported as an Unknown, Warning, Error or Fatal error status. The tool can also report an Operational error, which means there was an issue during the cluster health check. When the command encounters a Fatal or Operational error, it stops further diagnostics. For example, an inconsistency in the store metadata does not allow correct cluster operations and must be handled first.

Outputting the List of Unresolved Distributed Transactions

To view the list of distributed transactions that Postgres Pro Shardman built-in monitoring tools failed to resolve, run the following command:

```
shardmanctl [common_options] status transactions -r|--
repgroup replication_group_name
```

Each output transaction consists of `tx_id` (transaction ID), `coordinator_id`, `creation_time` and `description` (error or transaction status). To display the list of transactions for a specific replication group, use the `-r|--repgroup` option (for all replication groups by default). In case there are no such transactions, returns `null` value in JSON.

Store-related Commands

store dump

Syntax:

```
shardmanctl [common_options] store dump [-f|--file filename]
```

Dumps current cluster specifications from the store.

```
-f filename
--file=filename
```

Specifies the output file (– for stdout).

Default: –.

For more details, see [the section called “Dumping All Keys from the Store to Debug Error Configuration”](#).

store restore

Syntax:


```
shardmanctl [common_options] store restore [-f|--file filename] [-y|--yes]
```

Allows to safely restore the `etcd` cluster from the dump created with the `shardmanctl store dump`.

```
-f filename  
--file=filename
```

Specifies the name of the `etcd` keys dump.

```
-y  
--yes
```

Perform automatic confirmation.

store lock

Syntax:

```
shardmanctl [common_options] store lock [-f|--format text|json]
```

Shows the current cluster meta lock information.

```
-f=text|json  
--format=text|json
```

Specifies the output format.

Default: `text`.

For more details, see [the section called “Output Current Cluster Meta Lock Information”](#).

store get

Syntax:

```
shardmanctl [common_options] store get [[-a|--alias aliasname]|[-k|--key keyname]]  
[-f|--file filename]
```

Gets the specified key from the store.

```
-a aliasname  
--alias=ladle|cluster |spec|shardspec
```

Specifies the use of alias instead of the full key name. Cannot be used with `--key`.

For more details, see [the section called “Getting the Current Shard Specification”](#).

```
-k keyname  
--key=keyname
```

Specifies the key to retrieve from the store. Cannot be used with `--alias`.

For more details, see [the section called “Getting the Current Shard Specification”](#).

```
-f filename  
--file=filename
```

Specifies the file to print the value to.

Default: `-` (stdout).

For more details, see [the section called “Getting the Current Shard Specification”](#).

store keys

Syntax:

```
shardmanctl [common_options] store keys
```

Gets all keys with the current cluster prefix from the store.

For more details, see [the section called “Getting the Cluster and Ladle Key Names For the Current Cluster”](#).

store set

Syntax:

```
shardmanctl [common_options] store set [[-a|--alias aliasname]|[-k|--key keyname]]  
[-f|--file filename]
```

Creates or rewrites a key in the store.

```
-a ladle|cluster |spec|shardspec  
--alias=ladle|cluster |spec|shardspec
```

Specifies the use of alias instead of the full key name. Cannot be used with `--key`.

```
-k keyname  
--key=keyname
```

Specifies the key name to set in the store. Cannot be used with `--alias`.

```
-f filename  
--file=filename
```

Specifies the file with input data (`-` for stdin).

For more details, see [the section called “Setting a New Spec for the Cluster”](#).

Dumping All Keys from the Store to Debug Error Configuration

After facing an error while using Postgres Pro Shardman cluster, to fill in an exhaustive report, it is convenient to dump all specifications that could produce such an error with the following command:

```
shardmanctl [common_options] store dump -f|--file filename
```

Some harmless errors may be shown, but they will not interrupt dumping. If you do not specify the filename, dump will be sent to stdout and may pollute your terminal.

Getting the Current Shard Specification

To get the current shard specification, which is normally a part of cluster key in the store, use the following command:

```
shardmanctl [common_options] store get -a|--alias shardspec -f|--file filename
```

If the cluster key is corrupted itself, the shard specification will not be shown either. Instead of using the alias, you may also find out the full cluster data key name (by listing all keys with `store keys` command), use `store get` to retrieve it and find the `shard` part there. Mind that while using the last option, `shardman.config_uuid` parameter will not be deleted, which may result in a conflict in later use of this data; for manipulation with shard specification, it is recommended to use `shardmanctl store get -a shardspec` command.

Getting the Cluster and Ladle Key Names For the Current Cluster

To get all key names in the store at once, run the following command:

```
shardmanctl [common_options] store keys
```

It can only be shown in JSON format. It will also print alias names for keys that have them (excluding `shardspec` and `spec`, since they are parts of other keys)

Output Current Cluster Meta Lock Information

You can view information about current cluster meta locks that acquired by any command:

```
shardmanctl [common_options] store lock -f|--format json
```

To get the report in JSON format, use `-f|--format=json` option (omitted above since `text` format is used by default). In case the lock does not exists returns `Lock not found`

Setting a New Spec for the Cluster

To set a new spec part of the cluster specification, run the following command:

```
shardmanctl [common_options] store set --alias=spec --file=spec.json
```

Since `spec` is a part of cluster data key, it cannot be set with `--key`. If the provided file is not a valid JSON, the new spec part will not be set.

Table-related Commands

tables sharded info

Syntax:

```
shardmanctl [common_options] tables sharded info [-t|--table table_name]
```

Gets information about a sharded table.

```
-t table  
--table=table
```

Specifies the name of the table in the format `schema.table`

tables sharded list

Syntax:

```
shardmanctl [common_options] tables sharded list
```

Gets the list of all sharded tables.

tables sharded norebalance

Syntax:

```
shardmanctl [common_options] tables sharded norebalance
```

Gets the list of sharded tables with automatic rebalancing disabled.

tables sharded partmove

Syntax:

```
shardmanctl [common_options] tables sharded partmove [-t|--table table_name] [-s|--shard shard_name] [-p|--partnum number]
```

Moves the specified partition of a sharded table to a new shard.

```
-t table  
--table=table
```

Specifies the name of the table in the format `schema.table`.

```
-p number  
--partnum=number
```

Specifies the number of the partition to move.

```
-s shard_name  
--shard=shard_name
```

Specifies the name of the new shard for the partition.

tables sharded rebalance

Syntax:

```
shardmanctl [common_options] tables sharded rebalance [-t|--table table_name]
```

Enables and runs automatic data rebalancing for the selected sharded table.

```
-t table  
--table=table
```

Specifies the name of the table in the format `schema.table`.

Bench-related Commands

bench init

Syntax:

```
shardmanctl [common_options] bench init [--schema-type single|simple|shardman|custom]  
[--schema-file file_name] [-s|--scale scale_value] [-n|--no-vacuum]  
[-F|--fillfactor fillfactor_value]
```

Initializes the benchmark schema via pgbench. Schema can be custom or predefined. Creates `tpc-b` schema tables and fills them.

```
--schema-type=single|simple|shardman|custom
```

Type of schema used by schema initialization. Possible values:

- `single` — schema for a single PostgreSQL benchmark test
- `simple` — simple sharded schema
- `shardman` — sharded schema optimized for Postgres Pro Shardman
- `custom` — schema initialized by the user from the `--schema-file file`

Default schema: `shardman`.

```
--schema-file=file_name
```

File with DDL query for the custom schema type, to be used to create `tpc-b` tables for pgbench: `pgbench_accounts`, `pgbench_branches`, `pgbench_tellers`, `pgbench_history`.

```
-s scale_value
--scale=scale_value
```

Multiply the number of generated rows by the given scale factor.

```
-n
--no-vacuum
```

Perform no vacuuming during initialization.

```
-F fillfactor_value
--fillfactor=fillfactor_value
```

Fill pgbench tables with the given fillfactor value.

bench run

Syntax:

```
shardmanctl [common_options] bench run [--schema-type single|simple|shardman|custom]
[-f|--file file_name] [-c|--client client_value] [-C|--connect] [--full-output]
[-j|--jobs jobs_value] [-T|--time seconds] [-t|--transactions transactions_value]
[-s|--scale scale_factor] [-P | --progress seconds] [-R | --rate rate] [-M | --
protocol querymode]
```

Runs the initialized benchmark via pgbench. Can use the default pgbench script or a custom script from a file.

```
--schema-type=single|simple|shardman|custom
```

Type of schema used by schema initialization (bench init). Possible values:

- *single* — schema for single PostgreSQL benchmark
- *simple* — simple sharded schema
- *shardman* — sharded schema optimized for Postgres Pro Shardman
- *custom* — schema initialized by the user from the `--schema-file` file.

Default schema: *shardman*.

```
-f file_name
--file=file_name
```

Add a transaction script read from *filename* to the list of scripts to be executed.

Optionally, write an integer weight after @ to adjust the probability of selecting this script versus other ones. The default weight is 1. (To use a script file name that includes an @ character, append a weight so that there is no ambiguity, for example *filen@me@1*).

```
-c client_value
--client=client_value
```

Number of clients simulated, that is, number of concurrent database sessions.

```
-C
--connect
```

Establish a new connection for each transaction rather than doing it just once per client session.

```
--full-output
```

Print all pgbench output.

```
-j jobs_value
--jobs=jobs_value
```

Number of worker threads within pgbench.

`-s scale_factor`
`--scale=scale_factor`

Multiply the number of generated rows by the given scale factor.

`-T seconds`
`--time=seconds`

Run the test with a time limit of the specified number of seconds. The default value is 0, which means that there is no time limit. Note that the `-T` and `-t` flags are mutually exclusive, and the `-T` flag takes precedence. If the `-T` flag is specified, the `-t` flag is ignored.

`-t transactions_value`
`--transactions=transactions_value`

Run the test with a limit of the specified number of transactions on each client. The default value is 10. Note that the `-T` and `-t` flags are mutually exclusive, and the `-T` flag takes precedence. If the `-T` flag is specified, the `-t` flag is ignored.

`-P seconds`
`--progress=seconds`

Show progress report every *sec* seconds. The report includes the time since the beginning of the run, the TPS since the last report, and the transaction latency average, standard deviation, and the number of failed transactions since the last report. Under throttling (`-R`), the latency is computed with respect to the transaction scheduled start time, not the actual transaction beginning time, thus it also includes the average schedule lag time. When `--max-tries` is used to enable transaction retries after serialization/deadlock errors, the report includes the number of retried transactions and the sum of all retries.

`-R rate`
`--rate=rate`

Execute transactions targeting the specified rate instead of running as fast as possible (the default). The rate is given in transactions per second. If the targeted rate is above the maximum possible rate, the rate limit won't impact the results.

`-M querymode`
`--protocol=querymode`

Protocol to use for submitting queries to the server:

- `simple`: use simple query protocol.
- `extended`: use extended query protocol.
- `prepared`: use extended query protocol with prepared statements.

In the `prepared` mode, `pgbench` reuses the parse analysis result starting from the second query iteration, so `pgbench` runs faster than in other modes.

Default: `simple`.

bench cleanup

Syntax:

```
shardmanctl [common_options] bench cleanup
```

Cleans up schema database after benchmarks. Drops `tpc-b` tables.

bench generate

Syntax:

```
shardmanctl [common_options] bench generate [-c|--config file_name] [-o|--output-file file_name]
```

Gets the benchmark configuration from a file and generates a bash script to create a schema optimized for Postgres Pro Shardman and run the benchmark using pgbench. The configuration file must be in yaml format.

`-f file_name`
`--file=file_name`

The configuration file path. The file contains a sequence of script configurations. Each script must have a `schema_type: single|simple|shardman|custom`. For a custom schema it is necessary to specify the `schema_file` with the DDL script. Optional parameters: `init_flags` (default set: `-s 1000`), `run_flags` (default set: `-n -P 10 -c 10 -j 4 -T 60`), `partitions` (default value: 50). It is highly recommended to use `-n` (`--no-vacuum`) parameter inside `run_flags`. Configuration file example:

```
benches:
- schema_type: single
  init_flags: "-s 3"
  run_flags: "-n -P 10 -c 10 -j 4 -T 10"
- schema_type: simple
  init_flags: "-s 4"
  run_flags: "-n -P 10 -c 20 -j 4 -T 10"
  partitions: 100
- schema_type: shardman
  init_flags: "-s 5"
  run_flags: "-n -P 10 -c 20 -j 4 -T 10"
- schema_type: custom
  init_flags: "-s 6"
  schema_file: "schema.psql"
```

`-o file_name`
`--output-file=file_name`

Output file. Default: stdout.

Cleanup-related Commands

cleanup

Syntax:

```
shardmanctl [common_options] cleanup [-p|--processrepgroups] --after-node-operation|--after-rebalance
```

Performs cleanup after the nodes `add` or `rebalance` command.

`-p node_names`
`--processrepgroups=node_names`

Perform an actual cleanup. By default, the tool only shows actions to be done during the actual cleanup. For more details, see [the section called “Performing Cleanup”](#).

`--after-node-operation`

Perform cleanup after a failure of a nodes `add` command.

`--after-rebalance`

Perform cleanup after a failure of a `rebalance` command.

Performing Cleanup

By default, `cleanup` operates in the report-only mode, that is, the following command will only show actions to be done during actual cleanup:

```
shardmanctl [common_options] cleanup --after-node-operation|--after-rebalance
```

To perform the actual cleanup, run the following command:

```
shardmanctl [common_options] cleanup -p|--processrepgroups --after-node-operation|--after-rebalance
```

Schema-related Commands

schema verify

Syntax:

```
shardmanctl [common_options] schema verify [--filter partitions|extension|tables|roles|schemas|sequences]
```

Verifies that a schema has the correct information about partitions, and that they are located on the same primary nodes as the corresponding partitions of a colocated table. It also checks that all the shards have the same global and sharded tables, global roles, and global sequences, and that their names differ from the local versions, as well as checks that the shards have the same schemas and extensions.

```
--filter partitions|extension|tables|roles|schemas|sequences
```

Specifies what is included in the check. If not specified, all the checks are done.

Default: not specified.

schema list

Syntax:

```
shardmanctl [common_options] schema list [-f|--format text|json]
```

Outputs a list of the stored dump IDs of the `shardman` schema.

```
-f=text|json  
--format=text|json
```

Specifies the output format.

Default: text.

schema show

Syntax:

```
shardmanctl [common_options] schema show --shard shard_name --dump-id dump_id
```

Outputs a dump with the specified ID for the specified shard.

```
--shard=shard_name
```

Specifies the name of the shard.

```
--dump-id=dump_id
```

Dump ID.

schema restore

Syntax:


```
shardmanctl [common_options] schema restore --dump-id dump_id [--shards comma_separated_list_of_shards] [--smart]
```

Restores of a `shardman` schema from a dump with a specified ID.

`--dump-id=dump_id`

Dump ID.

`--shards=comma_separated_list_of_shards`

Specifies the list of the shards to restore. If not specified, all the cluster shards are restored. The dump must contain all shards that must be restored.

`--smart`

Restores the `shardman` schema from PostgreSQL system tables if the schema's integrity was violated. It can be used when the `shardmanSchemaDumpMonitorDisabled` configuration parameter is enabled, or if no valid dump for the schema is present. However, in some cases (e.g. if data was lost from all shards) no full functioning restoration is guaranteed, and it is advised to use the partly restored schema as data source to proceed with restoration to a different working Postgres Pro Shardman cluster. It requires the tables that are being restored to be empty, thus making the partial tables restore impossible.

schema dump

Syntax:

```
shardmanctl [common_options] schema dump [--no-verify]
```

Creates a `shardman` schema dump that is saved in the `etcd` with the `shardman/cluster_name/shardman_schema_dumps/dump_id` key. Maximum 10 dumps can be created. Should the limit be reached, the oldest dump is deleted. If a new dump duplicates one of the saved dumps, only the newest one remains.

`--no-verify`

Skip the consistency check of a dump.

Default: false (check is not skipped).

Daemon-related Commands

daemon set

Syntax:

```
shardmanctl [common_options] daemon set [--session-log-level debug | info | warn | error] [--session-log-format json|text] [--session-log-nodes]
```

Allows updating the log parameters “on the fly”.

`--session-log-level debug | info | warn | error`

Updates the log level to debug, info, warn, or error.

`--session-log-format json|text`

Updates the log output format to text or json.

`--session-log-nodes`

Specifies which cluster nodes must be updated. If not specified, the parameters are updated on every node.

Default: all nodes.

daemon check

Syntax:

```
shardmanctl [common_options] daemon check -n|--nodes node_name:port
```

Checks shardmand on nodes.

```
-n node_name:port
--nodes=node_name:port
```

List of nodes to check shardmand on. For more details, see [the section called “Checking shardmand Service on Nodes”](#).

Checking shardmand Service on Nodes

`daemon check` not only checks that shardmand service is running on specified nodes, but also assures those services are configured for the same cluster as shardmanctl:

```
shardmanctl [common_options] daemon check -n|--nodes node_names
```

Other Commands

getconnstr

Syntax:

```
shardmanctl [common_options] getconnstr --all
```

Gets the libpq connection string for connecting to a cluster as administrator.

```
--all
```

Adds replicas to `getconnstr`.

forall

Syntax:

```
shardmanctl [common_options] forall --sql query[ --sql query[ --sql query ...]] [--twophase] [-d --dbname]
```

Executes an SQL statement on all replication groups in a Postgres Pro Shardman cluster.

```
--sql query
```

Specifies the statement to be executed.

```
--twophase
```

Use the two-phase-commit protocol to execute the statement.

```
-d
```

```
--dbname
```

Specifies a database to run the request for. If not specified, the default value is used.

Default: postgres.

rebalance

Syntax:

```
shardmanctl [common_options] rebalance [-f|--force]
```

Rebalances sharded tables. Once the command is executed, a dump is created.

```
-f
--force
```

Perform forced rebalance of sharded tables whose partitions were manually moved.

upgrade

Syntax:

```
shardmanctl [common_options] upgrade
```

Upgrades the shardman database extension and updates `pg_foreign_server` options.

script

Syntax:

```
shardmanctl [common_options] script -s|--shard shard_name [[-f|--file file_name] [--sql query]]
```

Executes non-transactional commands from a file or from the command-line on the specified shards.

```
-s shard_name
--shard=shard_name
```

Shard name.

```
-f file_name
--file=file_name
```

Add a transaction script from the *file_name* file to the list of scripts to be executed.

```
--sql query
```

Specifies the statement to be executed and can only be used separately from `-f`.

psql

Syntax:

```
shardmanctl [common_options] psql -s|--shard shard_name
```

Connects to the first available primary node if no options are specified.

```
-s shard_name
--shard=shard_name
```

Name of the shard. If specified, the connection is installed with this shard current primary.

history

Syntax:

```
shardmanctl [common_options] history [--reverse | -r] [-f|--format json|text] [-l|--limit number_of_commands]
```

Shows history of the commands that updated the cluster. By default, they are sorted from the most recent to the oldest ones.

```
-r
--reverse
```

Switches to the ascending sorting order.

`-f json|text`
`--format=json|text`

Output format.

Default: `text`.

`-l`
`--limit=number_of_commands`

Limit for the number of the most recent commands in the output. The maximum value is 200.

Default: 20.

`--log-format json|text`

Specifies the log output format, `text` or `json`.

Common Options

`shardmanctl` common options are optional parameters that are not specific to the utility. They specify `etcd` connection settings, cluster name and a few more settings. By default `shardmanctl` tries to connect to the `etcd` store `127.0.0.1:2379` and use the `cluster0` cluster name. The default log level is `info`.

`-h, --help`

Show brief usage information.

`--cluster-name cluster_name`

Specifies the name for a cluster to operate on. The default is `cluster0`.

`--log-level level`

Specifies the log verbosity. Possible values of `level` are (from minimum to maximum): `error`, `warn`, `info` and `debug`. The default is `info`.

`--retries number`

Specifies how many times `shardmanctl` retries a failing `etcd` request. If an `etcd` request fails, most likely, due to a connectivity issue, `shardmanctl` retries it the specified number of times before reporting an error. The default is 5.

`--session-timeout seconds`

Specifies the session timeout for `shardmanctl` locks. If there is no connectivity between `shardmanctl` and the `etcd` store for the specified number of seconds, the lock is released. The default is 30.

`--store-endpoints string`

Specifies the `etcd` address in the format: `http[s]://address[:port](,http[s]://address[:port])*`. The default is `http://127.0.0.1:2379`.

`--store-ca-file string`

Verify the certificate of the HTTPS-enabled `etcd` store server using this CA bundle.

`--store-cert-file string`

Specifies the certificate file for client identification by the `etcd` store.

`--store-key string`

Specifies the private key file for client identification by the `etcd` store.

`--store-timeout duration`

Specifies the timeout for a `etcd` request. The default is 5 seconds.

`--monitor-port number`

Specifies the port for the shardmand http server for metrics and probes. The default is 15432.

`--api-port number`

Specifies the port for the shardmand http api server. The default is 15432.

`--version`

Show shardman-utils version information.

Environment

`SDM_BACKUP_MODE`

An alternative to setting the `--backup-mode` option.

`SDM_BACKUP_PATH`

An alternative to setting the `--backup-path` option.

`SDM_CLUSTER_NAME`

An alternative to setting the `--cluster-name` option.

`SDM_ETCD_PATH`

An alternative to setting the `--etcd-path` option.

`SDM_FILE`

An alternative to setting the `--file` option for config update.

`SDM_LOG_LEVEL`

An alternative to setting the `--log-level` option.

`SDM_NODES`

An alternative to setting the `--nodes` option for `nodes add` and `nodes rm`.

`SDM_RETRIES`

An alternative to setting the `--retries` option.

`SDM_SPEC_FILE`

An alternative to setting the `--spec-file` option for `init`.

`SDM_STORE_ENDPOINTS`

An alternative to setting the `--store-endpoints` option.

`SDM_STORE_CA_FILE`

An alternative to setting the `--store-ca-file` option.

`SDM_STORE_CERT_FILE`

An alternative to setting the `--store-cert-file` option.

`SDM_STORE_KEY`

An alternative to setting the `--store-key` option.

`SDM_STORE_TIMEOUT`

An alternative to setting the `--store-timeout` option.

SDM_SESSION_TIMEOUT

An alternative to setting the `--session-timeout` option.

Examples

Initializing the Cluster

To initialize a Postgres Pro Shardman cluster that has the `cluster0` name, uses an etcd cluster consisting of `n1,n2` and `n3` nodes listening on port 2379, ensure proper settings in the spec file `sdmspec.json` and run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 init -f sdmspec.json
```

Getting the Cluster Connection String

To get the connection string for a Postgres Pro Shardman cluster that has the `cluster0` name, uses an etcd cluster consisting of `n1,n2` and `n3` nodes listening on port 2379, run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 getconnstr
```

```
dbname=postgres host=n1,n4,n2,n1,n1,n2,n4,n3 password=yourpasswordhere
port=5432,5433,5432,5433,5432,5433,5432,5433 user=postgres
```

To add replicas to `getconnstr`, use `--all`.

Getting the Cluster Status

Here is a sample status output from `shardmanctl` with OK and Error statuses:

```
$ shardmanctl status --filter store,shardmand,rg --sort=node
```

== STORE STATUS ==			
STATUS NODE	MESSAGE	REPLICATION GROUP	
OK	etcd store is OK		
== SHARDMAND STATUS ==			
STATUS NODE	MESSAGE	REPLICATION GROUP	
OK 56d819b4e9e4	shardmand on node 56d819b4e9e4 is OK		
OK 6d0aabd50acc	shardmand on node 6d0aabd50acc is OK		
== REPLICATION GROUP STATUS ==			

STATUS NODE	MESSAGE	REPLICATION GROUP	
OK	Replication group clover-1-56d819b4e9e4 is	clover-1-56d819b4e9e4	
	OK		
	Replication connection is down for slave		
Error 6d0aabd50acc:5442	6d0aabd50acc:5442 in replication group clover-1-6d0aabd50acc	clover-1-6d0aabd50acc	
== RESTART REQUIRED PARAMS STATUS ==			
STATUS	MESSAGE	REPLICATION GROUP	NODE
OK	No pending restart parameters	shard-1	shrn1
OK	No pending restart parameters	shard-2	shrn4

Rewriting Shard Specification

First, get the list of available keys in the store using the following command:

```
$ shardmanctl store keys
```

```
$ shardmanctl store keys
```

```
{
  "Key": "shardman/cluster0/biha/cluster/shard-1/clusterdata"
}{
  "Key": "shardman/cluster0/biha/cluster/shard-1/keepers/info/keeper_1"
}{
  "Key": "shardman/cluster0/biha/cluster/shard-2/clusterdata"
}{
  "Key": "shardman/cluster0/biha/cluster/shard-2/keepers/info/keeper_1"
}{
  "Key": "shardman/cluster0/data/cluster",
  "Alias": "cluster"
}{
  "Key": "shardman/cluster0/data/cluster/revisions"
}{
  "Key": "shardman/cluster0/data/cluster/revisions/2024-10-31T09:29:26"
}{
  "Key": "shardman/cluster0/data/ladle",
  "Alias": "ladle"
}{
}
```

```
"Key": "shardman/cluster0/data/repgroups",
"Alias": "repgroups"
}
```

Get the shard configuration from the store and save it in the `shardspec.json` file with the command

```
$ shardmanctl store get -a shardspec -f shardspec.json
```

Apply the necessary changes to the file and upload the new specification using `shardmanctl config update`. **Mind that `shardman.config_uuid` parameter is deleted with `shardmanctl store get -a shardspec` and not with `shardmanctl store get -k full/path/to/clusterspec`; using spec with existing `shardman.config_uuid` will result in a conflict.**

Important

Do not use `store set` command to update cluster configurations because it *does not* apply a new specification on all nodes, it only writes it to the store. For the above example with shard specification, `shardmanctl config update` or `shardmanctl set` is acceptable.

To double-check, you can get the cluster key with new `StolonSpec` by the full key name (which was shown earlier with `store keys` command):

```
$ shardmanctl store get -k shardman/cluster0/data/cluster
```

```
{
  "FormatVersion": 1,
  "Spec": {
    "PgSuAuthMethod": "md5",
    "PgSuPassword": "12345",
    "PgSuUsername": "postgres",
    "PgReplAuthMethod": "md5",
    "PgReplPassword": "12345",
    "PgReplUsername": "repluser",
    "ShardSpec": {
      ...
    }
  }
}
```

Adding Nodes to the Cluster

To add `n1,n2,n3` and `n4` nodes to the cluster, run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 nodes add
-n n1,n2,n3,n4
```

The `n1,n2,n3` and `n4` nodes are added as primaries for the replication groups. After that you can get the cluster topology:

```
$ shardmanctl cluster topology
```

```
== REPLICATION GROUP shard-1, RGID - 1 ==
+-----+-----+-----+-----+
| HOST   | PORT | KEEPER | STATUS |
+-----+-----+-----+-----+
| n1     | 2379 | yes    | primary |
+-----+-----+-----+-----+
| n2     | 2379 | yes    | primary |
+-----+-----+-----+-----+
| n3     | 2379 | yes    | primary |
+-----+-----+-----+-----+
| n4     | 2379 | yes    | primary |
+-----+-----+-----+-----+
| PRIORITY |
+-----+-----+-----+-----+
```


shardmanctl

n1	5432	keeper_1	LEADER_RW
NOT SET			
== REPLICATION GROUP shard-2, RGID - 2 ==			
HOST	PORT	KEEPER	STATUS
PRIORITY			
n2	5432	keeper_1	LEADER_RW
NOT SET			
== REPLICATION GROUP shard-3, RGID - 3 ==			
HOST	PORT	KEEPER	STATUS
PRIORITY			
n3	5432	keeper_1	LEADER_RW
NOT SET			
== REPLICATION GROUP shard-4, RGID - 4 ==			
HOST	PORT	KEEPER	STATUS
PRIORITY			
n4	5432	keeper_1	LEADER_RW
NOT SET			

To add a new replica, run:

```
shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 shard -s  
shard-1 add -n n2
```

This will give a cluster with the following topology:

```
$ shardmanctl cluster topology
```

== REPLICATION GROUP shard-1, RGID - 1 ==			
HOST	PORT	KEEPER	STATUS
PRIORITY			
n1	5432	keeper_1	LEADER_RW
NOT SET			
n2	5433	keeper_2	FOLLOWER
NOT SET			
== REPLICATION GROUP shard-2, RGID - 2 ==			

HOST PRIORITY	PORT	KEEPER	STATUS	
n2 NOT SET	5432	keeper_1	LEADER_RW	
== REPLICATION GROUP shard-3, RGID - 3 ==				
HOST PRIORITY	PORT	KEEPER	STATUS	
n3 NOT SET	5432	keeper_1	LEADER_RW	
== REPLICATION GROUP shard-4, RGID - 4 ==				
HOST PRIORITY	PORT	KEEPER	STATUS	
n4 NOT SET	5432	keeper_1	LEADER_RW	

Removing Nodes from the Cluster

To remove `n1` and `n2` nodes from the `cluster0` cluster, run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379
nodes rm -n n1,n2
```

If `cross` placement policy is used, then the clovers that contain them will be deleted along with the nodes.

Executing a Query on All Replication Groups

To execute the `select version()` query on all replication groups, run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 forall --
sql 'select version()'
```

Node 1 says:

```
[PostgreSQL 13.1 on x86_64-pc-linux-gnu, compiled by gcc (Ubuntu
9.3.0-17ubuntu1~20.04) 9.3.0, 64-bit]
```

Node 4 says:

```
[PostgreSQL 13.1 on x86_64-pc-linux-gnu, compiled by gcc (Ubuntu
9.3.0-17ubuntu1~20.04) 9.3.0, 64-bit]
```

Node 3 says:

```
[PostgreSQL 13.1 on x86_64-pc-linux-gnu, compiled by gcc (Ubuntu
9.3.0-17ubuntu1~20.04) 9.3.0, 64-bit]
```

Node 2 says:

```
[PostgreSQL 13.1 on x86_64-pc-linux-gnu, compiled by gcc (Ubuntu
9.3.0-17ubuntu1~20.04) 9.3.0, 64-bit]
```

Performing Rebalance

To rebalance sharded tables in the `cluster0` cluster, run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 rebalance
```

Updating PostgreSQL Configuration Settings

To set the `max_connections` parameter to 200 in the cluster, create the spec file (for instance, `~/shard-spec.json`) with the following contents:

```
{
  "pgParameters": {
    "max_connections": "200"
  }
}
```

Then run:

```
$ shardmanctl --store-endpoints http://n1:2379,http://n2:2379,http://n3:2379 config
update -p -f ~/shardspec.json
```

Since changing `max_connections` requires a restart, DBMS instances are restarted by this command.

Performing Backup and Recovery

To create a backup of the `cluster0` cluster using `etcd` at `etcdserver` listening on port 2379 and store it in the local directory `/var/backup/shardman`, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 backup --datadir=/var/backup/
shardman --use-ssh
```

Assume that you are performing a recovery from a backup to the `cluster0` cluster using `etcd` at `etcdserver` listening on port 2379 and you take the backup description from the `/var/backup/shardman/backup_info` file. Edit the `/var/backup/shardman/backup_info` file, set `DataRestoreCommand`, `RestoreCommand` as necessary and run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 recover --info /var/backup/
shardman/backup_info
```

For metadata-only restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 recover --metadata-only --
dumpfile /var/backup/shardman/etcd_dump
```

For schema-only restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 recover --schema-only --
dumpfile /var/backup/shardman/etcd_dump
```

For single shard restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 recover --info /var/backup/
shardman/backup_info --shard shard_1
```

Performing Backup and Recovery with probackup Command

To create a backup of the `cluster0` cluster using `etcd` at `etcdserver` listening on port 2379 and store it in the local directory `/var/backup/shardman`, first initialize the backups repository with the `init` subcommand:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup init --backup-path=/
var/backup/shardman --etcd-path=/var/backup/etcd_dump
```

Then add and enable `archive_command` with the `archive-command` subcommand:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup archive-command add --
backup-path=/var/backup/shardman
```

If the repository is successfully initialized and `archive-command` successfully added, create a `FULL` backup with the `backup` subcommand:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup backup --backup-path=/
var/backup/shardman --etcd-path=/var/backup/etcd_dump --backup-mode=FULL --compress --
compress-algorithm=zlib --compress-level=5
```

To create `DELTA`, `PTRACK` or `PAGE` backup, run the `backup` subcommand with `DELTA`, `PTRACK` or `PAGE` value of the `--backup-mode` option:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup backup --backup-path=/
var/backup/shardman --etcd-path=/var/backup/etcd_dump --backup-mode=DELTA --compress --
compress-algorithm=zlib --compress-level=5
```

To show the created backup ID, run `show` subcommand:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup show --backup-path=/
var/backup/shardman --format=table
```

== BACKUP ID 'S88FRO'									
==									
INSTANCE				HOST					
RECOVERY TIME		MODE		WAL MODE		TLI		DATA	
Z-RATIO		START LSN		STOP LSN		STATUS		WAL	
shard-1				n1					
2024-02-02 14:19:05+00		FULL		ARCHIVE		1/0		42.37MiB	
16MiB		1.00		0/C000028		0/D0018B0		OK	
shard-2				n2					
2024-02-02 14:19:05+00		FULL		ARCHIVE		1/0		42.38MiB	
16MiB		1.00		0/C000028		0/D001E00		OK	

In `PTRACK` backup mode, Postgres Pro Shardman tracks page changes on the fly. Continuous archiving is not necessary for it to operate. Each time a relation page is updated, this page is marked in a special `PTRACK` bitmap. Tracking implies some minor overhead on the database server operation, but speeds up incremental backups significantly.

If you are going to use `PTRACK` backups, complete the following additional steps:

- Preload the `ptrack` shared library on each node. This can be done by adding the `ptrack` value to the `shared_preload_libraries` parameter.
- Create the `PTRACK` extension on each cluster node:

```
$ shardmanctl --store-endpoints http://etcdserver:2379
forall --sql "create extension ptrack"
```

- To enable tracking page updates, set the `ptrack.map_size` parameter as follows:

```
$ shardmanctl --store-endpoints http://etcdserver:2379
update '{"pgParameters":{"ptrack.map_size":"64"}}'
```

For optimal performance, it is recommended to set `ptrack.map_size` to $N/1024$, where N is the maximum size of the cluster node, in MB. If you set this parameter to a lower value, PTRACK is more likely to map several blocks together, which leads to false-positive results when tracking changed blocks and increases the incremental backup size as unchanged blocks can also be copied into the incremental backup. Setting `ptrack.map_size` to a higher value does not affect PTRACK operation, but it is not recommended to set this parameter to a value higher than 1024.

To validate the created backup, run `validate` subcommand:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup validate --backup-
path=/var/backup/shardman --backup-id=RFP1FI
```

Assume that you are performing a recovery from a backup to the `cluster0` cluster using `etcd` at `etcd-server` listening on port 2379 and you take the backup ID from the `show` command:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup restore --backup-
path=/var/backup/shardman --backup-id=RFP1FI
```

Finally we need to enable `archive_command` back.

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup archive-command add --
backup-path=/var/backup/shardman
```

For metadata-only restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup restore --metadata-
only --backup-path=/var/backup/shardman --backup-id=RFP1FI
```

For metadata-only restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup restore --schema-only
--backup-path=/var/backup/shardman --backup-id=RFP1FI
```

For single shard restore, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup restore --backup-
path=/var/backup/shardman --backup-id=RFP1FI --shard shard_1
```

For Point-in-Time Recovery, run:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 probackup restore --metadata-
only --backup-path=/var/backup/shardman --backup-id=RFP1FI --recovery-target-
time='2006-01-02 15:04:05' -s
```

Loading Data from a Text File

To load data into a Postgres Pro Shardman cluster, run the following command:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 load --file=/var/load/data.tsv
--table=mytable --source file --format text -j 8
```

In this example, data is loaded from the `/var/load/data.tsv` data file (tab-delimited) into the table `mytable` in 8 parallel threads. You can use `schema.table` as the table name.

Loading data from PostgreSQL table

To load data into a Postgres Pro Shardman cluster from a PostgreSQL table, run the following command:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 load -t desttable --
source postgres --source-connstr "dbname=db host=srchost port=srcport user=login
password=passwd" --source-table sourcetable -j 8
```

In this example, data is loaded from the table `sourcetable` into the `desttable` table in 8 parallel threads. You can use `schema.table` as table names.

Loading Data with a Schema from PostgreSQL

To load data with a schema into Postgres Pro Shardman cluster from PostgreSQL, run the following command:

```
$ shardmanctl --store-endpoints http://etcdserver:2379 load --schema load_schema.yaml
```

The file `load_schema.yaml` has the following format:

```
version: "1.0"
migrate:
  connstr: "dbname=workdb host=workhost port=workport user=workuser
password=workpassword"
  jobs: 8
  batch: 1000
  options:
    - create_schema
    - create_table
    - create_index
    - create_sequence
    - create_foreign_key
    - create_role
    - copy_ownership
    - copy_grants
    - truncate_table
    - skip_no_pkey_tables
    - skip_create_index_error
    - skip_create_extension_error
    - skip_load_errors
    - skip_create_foreign_key_error
    - skip_create_role_error
    - skip_copy_grants_error
    - skip_copy_ownership_error
  schemas:
    - name: public
      all: false
      tables:
        - name: tab1
          type: sharded
          partitions: 6
          distributedby: id
          priority: 3
        - name: tab2
          type: global
        - name: tab3
          type: sharded
          partitions: 6
          distributedby: field_id
          colocatewith: tab1
        - name: table4
          type: global
          source: schema.view
          source_pk: field_id
        - name: table5
          type: global
          source: schema.func(arg)
```

```
        source_pk: field_id
- name: schema2
  all: false
  default_type: sharded
  default_partitions: 6
  tables:
    - name: table1
      distributedby: field_id
      priority: 2
    - name: table2
      type: global
    - name: table3
      source: schema.view
      distributedby: field_id
      priority: 3
    - name: table4
      distributedby: field_id
      source: schema.func(arg)
    - name: table5
      source: schema."complex."table.name"
      distributedby: field_id
- name: schema3
  all: true
  skip_tables: [table1, table2, table3]
roles:
- name: test_user1
  password: test_password
- name: test_user2
```

The `migrate.jobs` value defines the number of parallel data loader processes.

The `migrate.batch` value is the number of rows in one batch (recommended value is 1000).

The `migrate.schemas` section defines an array of source database schemas that you are working with. All other schemas will be skipped.

If the `all` value is set to `true`, then all tables from the current schema will be migrated (with `global` type by default). If a table is listed in the `migrate.schemas.tables` array, then the target table type must be explicitly specified for it. Two types of tables are currently supported: `global` and `sharded`. Global tables are loaded first, then `sharded` tables and at the end `sharded` tables with the `colocatedwith` parameter. The order of loading tables of the same type can be changed using `priority` option.

The `migrate.schemas.skip_tables` section defines an array of table names that will be skipped when the schema is loaded even if the `all` parameter is set to `true`.

For `sharded` tables, the following attributes must be set: `distributedby` (specifies the name of the column to use for the table partitioning) and `partitions` (number of partitions that will be created for this table). Optionally, for `sharded` tables `colocatewith` attribute can be set (name of the table to colocate with). Postgres Pro Shardman will try to place partitions of the created table with the same partition key on the same nodes as the corresponding partitions of the table specified by `colocatewith`.

You can specify the table `default_type` option for a schema: `global` or `sharded` (default: `global`). For the `sharded` type you can also specify the `default_partitions` option (default: 20). If you set `default_type` to `sharded`, you need to specify the `distributedby` option for each table.

The `source` option for a table should include the schema and table source: `schema.source`. The source can be a table, view or function. For example: `public.table`, `public.view`, `public.func(arg)`. If you set the `source` view or function for a `global` table, you should specify `source_pk` to set the primary key

for this table. If `source` is not specified or contains the name of a table, you can also specify `source_pk` to create a primary key or override the existing one.

The `priority` option for table determines the order in which the tables of the same type are loaded. Tables with higher `priority` are loaded earlier. Default `priority` value is 0.

The `migrate.roles` section defines an array of role names and passwords that will be copied from the source database if `create_role` is specified.

The schema supports the following options:

- `create_schema` — create database schemas if they do not exist.
- `create_table` — create tables if they do not exist.
- `create_index` — create indexes after creating tables.
- `create_sequence` — create sequences if they do not exist.
- `create_foreign_key` — create foreign keys after creating tables.
- `truncate_table` — truncate tables before data load.
- `create_role` — create global roles defined in `migrate.roles` and copy role parameters from the source database.
- `copy_grants` — copy access privileges from the source database.
- `copy_ownership` — change of table owners to the owner in the source database.
- `skip_no_pkey_tables` — skip tables without primary keys.
- `skip_create_index_error` — skip index creation errors.
- `skip_create_extension_error` — skip extension creation errors.
- `skip_load_errors` — continue loading if errors occur.
- `skip_create_foreign_key_error` — skip foreign key creation errors.
- `skip_create_role_error` — skip role creation errors.
- `skip_copy_ownership_error` — skip table owner changing errors.
- `skip_copy_grants_error` — skip errors when copying access privileges from the source database.

Initialization and Running Benchmarks

To initialize a benchmark via `shardmanctl` using `pgbench` with the `shardman` schema, `scale=1000`, `partitions=40`, run:

```
$ shardmanctl bench init --schema-type=shardman --scale=1000 --partitions=40
```

To run an initialized benchmark for the same `shardman` schema, number of jobs=4, number of clients=10, duration in seconds=60 and full `pgbench` output, use:

```
$ shardmanctl bench run --schema-type=shardman --jobs=4 --client=10 --time=60 --full-output
```

To initialize a benchmark with the `custom` schema from file `schema.psql` with `scale=1000` run:

```
$ shardmanctl bench init --schema-type=custom --schema-file=schema.psql --scale=1000
```

To run an initialized benchmark with the `custom` schema and custom transaction script from `script.psql` with the number of jobs=4, number of clients=10, duration in seconds=60, use:

```
$ shardmanctl bench run --schema-type=custom --file=script.psql --jobs=4 --client=10 --time=60
```


To clean up a PostgreSQL database of `tpc-b` tables, use:

```
$ shardmanctl bench cleanup
```

Benchmark Generation Scripts

To generate a benchmark sequence via `shardmanctl` from the config file `cfg.yaml` and output the result to file `script.sh`, run:

```
$ shardmanctl bench generate --config=cfg.yaml --output-file=script.sh
```

Configuration file example:

```
benches:
- schema_type: single
  init_flags: "-s 3"
  run_flags: "-n -P 10 -c 10 -j 4 -T 10"
- schema_type: simple
  init_flags: "-s 4"
  run_flags: "-n -P 10 -c 20 -j 4 -T 10"
  partitions: 100
- schema_type: shardman
  init_flags: "-s 5"
  run_flags: "-n -P 10 -c 20 -j 4 -T 10"
- schema_type: custom
  init_flags: "-s 6"
  schema_file: "schema.psql"
```

shardmand

shardmand — Postgres Pro Shardman configuration daemon

Postgres Pro Shardman Start-up with shardmand

shardmand is a Postgres Pro Shardman configuration daemon. It runs on each node in a Shardman cluster, generates a `boot_uuid` at every startup, subscribes for changes of `shardman/cluster0/data/ladle` and `shardman/cluster0/data/cluster` keys in the etcd store (`cluster0` is the default cluster name used by Shardman utils), and manages Postgres Pro Shardman processes on the node where it is running according to the configuration described in these JSON documents.

shardmand manages integrated `keepers` and `sentinels`. On startup and when one of the monitored etcd keys changes, shardmand reconfigures them as follows:

- It calculates the expected node configuration, i. e., the list of `keepers` and `sentinels` expected to run and their configurations, from the `shardman/cluster0/data/ladle` and `shardman/cluster0/data/cluster` values.
- It receives the list of running `keepers` and `sentinels` with their configurations from the internal process manager.
- It stops processes that are not expected to run. This can be a process that belongs to a cluster with the same name, but a different UUID, or a process whose description is no longer present in the expected node configuration. For `keeper` processes, shardmand purges their data directory.
- If a process should be running, but its settings are different from the expected ones, shardmand updates the configuration and restarts the process. If a process should be running, but it is not running, shardmand starts it.

Also, a separate thread of shardmand periodically updates the `shardman/cluster0/data/shardmand/NODENAME` etcd key with the `ClusterUUID` of the last cluster to which the configuration was applied. So, before the [shardmanctl](#) nodes add command tries to initialize new clusters for a clover, the command can ensure that no alive threads from a previous cluster configuration are left on all nodes in the clover.

Additionally, shardmand starts two http servers in separate threads. If servers ports match, a single server running both roles is started. The first server provides following metrics: `shardmand_etcd_unavailable_time_seconds`, `shardmand_healthy_keepers`, `shardmand_sentinels`, `shardmand_uptime`, `shardmand_etcd_errors_total`, `shardmand_reconfigurations_number_total`, `shardmand_demotions_number_total`. Also server provides a `/healthz` endpoint for shardmand health-check. The second server provides the following endpoints:

- `/shardmand/v1/replica` — returns 200 status code if a secondary instance is running on node, 500 status code if a primary instance is running on node, `/shardmand/v1/master` — returns 200 status code if a primary instance is running on node, 500 status code if a secondary instance is running on node, `/shardmand/v1/referee` — returns 200 status code if the only instance running on node is referee, 500 status code if the only instance running on node is primary or secondary, 404 if more than one instance is running. If node both primary and secondary instances are running on node `/shardmand/v1/replica` and `shardmand/v1/master` endpoints return 404 status code.
- `/shardmand/v1/status` — getting information about shardmand status, including `boot_uuid`.

All Postgres Pro Shardman services are managed by `shardmand@cluster0.service`, so when it is started, stopped, or restarted, it also starts, stops, or restarts all other Postgres Pro Shardman processes (including DBMS instances).

shardmand Syntax

```
shardmand [common_options] [ --system-bus ] [ --user user_name ]
```

Here *common_options* are:

```
[ --cluster-name cluster_name ] [ --log-level error | warn | info | debug ] [ --retries retries_number ] [ --session-timeout seconds ] [ --store-endpoints store_endpoints ] [ --store-ca-file store_ca_file ] [ --store-cert-file store_cert_file ] [ --store-key client_private_key ] [ --store-timeout duration ] [ --version ] [ -h | --help ] [ --log-format ]
```

Command-line Reference

This refsection describes `shardmand`-specific command-line options.

`--log-format`

Specifies the log output format, `json` or `text`. The default is `text`.

`--system-bus`

Not used. Left for compatibility. Ignored.

`--user user_name`

Not used. Left for compatibility. Ignored.

Common Options

`shardmand` common options are optional parameters that are not specific to the utility. They specify `etcd` connection settings, cluster name and a few more settings. By default `shardmand` tries to connect to the `etcd` store `127.0.0.1:2379` and use the `cluster0` cluster name. The default log level is `info`.

`-h, --help`

Show brief usage information.

`--cluster-name cluster_name`

Specifies the name for a cluster to operate on. The default is `cluster0`.

`--log-level level`

Specifies the log verbosity. Possible values of *level* are (from minimum to maximum): `error`, `warn`, `info` and `debug`. The default is `info`.

`--retries number`

Specifies how many times `shardmanctl` retries a failing `etcd` request. If an `etcd` request fails, most likely, due to a connectivity issue, `shardmanctl` retries it the specified number of times before reporting an error. The default is 5.

`--session-timeout seconds`

Specifies the session timeout for `shardmanctl` locks. If there is no connectivity between `shardmanctl` and the `etcd` store for the specified number of seconds, the lock is released. The default is 30.

`--store-endpoints string`

Specifies the `etcd` address in the format: `http[s]://address[:port](,http[s]://address[:port])*`. The default is `http://127.0.0.1:2379`.

`--store-ca-file string`

Verify the certificate of the HTTPS-enabled `etcd` store server using this CA bundle.

`--store-cert-file string`

Specifies the certificate file for client identification by the `etcd` store.

`--store-key string`

Specifies the private key file for client identification by the `etcd` store.

`--store-timeout duration`

Specifies the timeout for a etcd request. The default is 5 seconds.

`--monitor-port number`

Specifies the port for the shardmand http server for metrics and probes. The default is 15432.

`--api-port number`

Specifies the port for the shardmand http api server. The default is 15432.

`--version`

Show shardman-utils version information.

shardmand Environment

A shardmand service reads the environment from `/etc/shardman/shardmand-cluster0.env`. The following environment variables affect the behavior of shardmand.

`SDM_CLUSTER_NAME`

An alternative to setting the `--cluster-name` option

`SDM_LOG_LEVEL`

An alternative to setting the `--log-level` option

`SDM_RETRIES`

An alternative to setting the `--retries` option

`SDM_SYSTEM_BUS`

An alternative to setting the `--system-bus` option

`SDM_STORE_ENDPOINTS`

An alternative to setting the `--store-endpoints` option

`SDM_STORE_CA_FILE`

An alternative to setting the `--store-ca-file` option

`SDM_STORE_CERT_FILE`

An alternative to setting the `--store-cert-file` option

`SDM_STORE_KEY`

An alternative to setting the `--store-key` option

`SDM_STORE_TIMEOUT`

An alternative to setting the `--store-timeout` option

`SDM_SESSION_TIMEOUT`

An alternative to setting the `--session-timeout` option

`SDM_USER`

An alternative to setting the `--user` option

Part VII. Internals

This part contains assorted information that might be of use to Postgres Pro developers.

Chapter 53. Overview of Postgres Pro Internals

Author

This chapter originated as part of [sim98](#) Stefan Simkovics' Master's Thesis prepared at Vienna University of Technology under the direction of O.Univ.Prof.Dr. Georg Gottlob and Univ.Ass. Mag. Katrin Seyr.

This chapter gives an overview of the internal structure of the backend of Postgres Pro. After having read the following sections you should have an idea of how a query is processed. This chapter is intended to help the reader understand the general sequence of operations that occur within the backend from the point at which a query is received, to the point at which the results are returned to the client.

53.1. The Path of a Query

Here we give a short overview of the stages a query has to pass to obtain a result.

1. A connection from an application program to the Postgres Pro server has to be established. The application program transmits a query to the server and waits to receive the results sent back by the server.
2. The *parser stage* checks the query transmitted by the application program for correct syntax and creates a *query tree*.
3. The *rewrite system* takes the query tree created by the parser stage and looks for any *rules* (stored in the *system catalogs*) to apply to the query tree. It performs the transformations given in the *rule bodies*.

One application of the rewrite system is in the realization of *views*. Whenever a query against a view (i.e., a *virtual table*) is made, the rewrite system rewrites the user's query to a query that accesses the *base tables* given in the *view definition* instead.

4. The *planner/optimizer* takes the (rewritten) query tree and creates a *query plan* that will be the input to the *executor*.

It does so by first creating all possible *paths* leading to the same result. For example if there is an index on a relation to be scanned, there are two paths for the scan. One possibility is a simple sequential scan and the other possibility is to use the index. Next the cost for the execution of each path is estimated and the cheapest path is chosen. The cheapest path is expanded into a complete plan that the executor can use.

5. The executor recursively steps through the *plan tree* and retrieves rows in the way represented by the plan. The executor makes use of the *storage system* while scanning relations, performs *sorts* and *joins*, evaluates *qualifications* and finally hands back the rows derived.

In the following sections we will cover each of the above listed items in more detail to give a better understanding of Postgres Pro's internal control and data structures.

53.2. Query Processing

Postgres Pro Shardman uses the standard PostgreSQL [Section 53.1](#). Other nodes in the cluster are accessed via the modified `postgres_fdw` extension.

Postgres Pro Shardman query planner takes the query abstract syntax tree (AST) and creates a query plan, which is used by the executor. While evaluating query execution methods, the planner operates with so-called paths, which specify how relations should be accessed. While processing a query join tree, the planner looks at different combinations of how relations can be joined. Each time it examines a join

of two relations, one of which can be a join relation itself. After choosing the order and strategies for joining relations the planner considers the group by, order by and limit operations. When the cheapest path is selected, it is transformed to a query plan. A plan consists of a tree of nodes, each of which has methods to get one next result row (or NULL if there are no more results).

53.2.1. Push-down Technique

53.2.1.1. Joins

The efficiency of query execution in a distributed DBMS is determined by how many operations can be executed on nodes that hold the actual data. For Postgres Pro Shardman, a lot of effort is devoted to pushing down join operations. When the planner finds a relation that is accessible via a foreign data wrapper (FDW), it creates `ForeignPath` to access it. Later, when it examines a join of two relations and both of them are available via `ForeignPath` from the same foreign server, it can consider pushing down this join to the server and generating a so-called `ForeignJoinPath`. The planner can fail to do it if the join type is not supported, if filters attached to the relation should be applied locally, or if the relation scan result contains fields that cannot be evaluated on the remote server. An example of a currently unsupported join type is right anti-join. Local filters attached to the relation should be applied locally when remote execution can lead to a different result or if the `postgres_fdw` module cannot create SQL expressions to apply some of the filters. An example of fields that cannot be evaluated on a remote server are attributes of semi-join inner relation that are not accessible via an outer relation. If the `foreign_join_fast_path` configuration parameter is set to on (which is the default value), the Postgres Pro Shardman planner stops searching for other join strategies of two relations once it finds a foreign join possible for them. When the `postgres_fdw.enforce_foreign_join` configuration parameter is set to on (which is also the default), the cost of a foreign join is estimated so as to be always less than the cost of a local join.

When several sharded tables are joined on a sharding key, a partitionwise join can be possible. This means that instead of joining original tables, we can join their matching partitions. Partitionwise join currently applies only when the join conditions include all the partition keys, which must be of the same data type and have exactly matching sets of child partitions. Partitionwise join is crucial to the efficient query execution as it allows pushing down joins of table partitions. Evidently, to push down a join of several partitions, these partitions should reside on the same node. This is usually the case when sharded tables are created with the same `num_parts` parameter. However, for a rebalance process to move the corresponding partitions to the same nodes, sharded tables should be marked as colocated when created (see [Section 70.1.1](#)). Partitionwise join is enabled with the `enable_partitionwise_join` configuration parameter, which is turned on by default in Postgres Pro Shardman.

When a sharded table is joined to a global table, asymmetric partitionwise join is possible. This means that instead of joining original tables, we can join each partition of the sharded table with the global table. This makes it possible to push down a join of sharded table partitions - with a global table to the foreign server.

53.2.1.2. Aggregations

After planning joins, the planner considers paths for post-join operations, such as aggregations, limiting, sorting and grouping. Not all such operations reach FDW pushdown logic. For example, pushing down the `LIMIT` clause to a remote server is not possible when additional sorting is necessary prior to merging results from different nodes. There are two efficient strategies for executing aggregates on remote nodes. The first one is a partitionwise aggregation — when a `GROUP BY` clause includes a partitioning key, the aggregate can be pushed down together with the `GROUP BY` clause (this behavior is controlled by the `enable_presorted_aggregate`/`enable_partitionwise_aggregate` configuration parameter, which is turned on by default in Postgres Pro Shardman). Alternatively, the planner can decide to execute partial aggregation on each partition of a sharded table and then combine the results. If necessary, the executor can export the representation of the aggregate's internal state on each node, import it locally and execute the final aggregation steps on coordinator. Usually it works for the built-in aggregates. Also the planner refuses pushing down partial aggregates if they contain additional clauses, such as `ORDER BY` or `DISTINCT`, or if the statement has the `HAVING` clause.

53.2.1.3. Subqueries

Generally, subqueries cannot be pushed down to other cluster nodes. However, Postgres Pro Shardman uses subquery unnesting to alleviate this limitation.

In PostgreSQL, non-correlated subqueries can be transformed into semi-joins. In the following example, ANY subquery on non-partitioned tables is transformed to Hash Semi Join:

```
EXPLAIN (COSTS OFF) SELECT * FROM pgbench_branches WHERE bid = ANY (SELECT bid FROM
pgbench_tellers);
```

QUERY PLAN

```
-----
Hash Semi Join
  Hash Cond: (pgbench_branches.bid = pgbench_tellers.bid)
    -> Seq Scan on pgbench_branches
    -> Hash
        -> Seq Scan on pgbench_tellers
```

Planner also can convert some forms of correlated subqueries (i.e., subqueries that reference upper-level relations) into semi-joins. This optimization works for IN and = operators. The transformation has some restrictions. For example, it is not considered if a subquery contains aggregates or references upper-level relations from outside of a WHERE clause. This optimization allows transforming more complex subqueries into semi-joins, like in the following example:

```
EXPLAIN (COSTS OFF) SELECT * FROM pgbench_branches WHERE bid = ANY (SELECT bid FROM
pgbench_tellers WHERE tbalance = bbalance);
```

QUERY PLAN

```
-----
Hash Semi Join
  Hash Cond: ((pgbench_branches.bid = pgbench_tellers.bid) AND
(pgbench_branches.bbalance = pgbench_tellers.tbalance))
    -> Seq Scan on pgbench_branches
    -> Hash
        -> Seq Scan on pgbench_tellers
(5 rows)
```

After applying subquery unnesting, semi-join can be pushed down for execution to a remote node.

53.2.2. Asynchronous Execution

When a sharded table is queried, the Postgres Pro Shardman planner creates `Append` plans to scan all partitions of the table and combine the result. When some of partitions are foreign tables, the planner can decide to use an asynchronous execution. This means that when an `Append` or `MergeAppend` node for the first time after initialization is asked for the tuples, it asks asynchronous child nodes to start fetching the result. For `postgres_fdw` `async ForeignScan` nodes, it means that a remote cursor is declared and a fetch request is sent to the remote server. If Silk transport is used, this means that the query is sent for execution to the remote server as an `MT_SPI` message.

After sending a request to the remote servers, `Append` returns to fetching data from synchronous child nodes — local scan nodes or synchronous `ForeignScan` nodes. Data from such nodes is fetched in a blocking manner. When `Append` ends getting data from synchronous nodes (such as local partitions scans), it looks if `async` nodes have some data. If they do not, it waits for `async` nodes to produce results.

Postgres Pro Shardman can execute several types of plans asynchronously. These are asynchronous `ForeignScans`, projections and trivial subquery scans (`select * from subquery`) over asynchronous plans.

The asynchronous execution is turned on by default on the level of a foreign server. This is controlled by `async_capable postgres_fdw` option. Both `Append` and `MergeAppend` plans support asynchronous execution.

Examples:

```
EXPLAIN (COSTS OFF) SELECT * FROM pgbench_accounts;  
QUERY PLAN
```

Append

```
-> Seq Scan on pgbench_accounts_0 pgbench_accounts_1  
-> Async Foreign Scan on pgbench_accounts_1_fdw pgbench_accounts_2  
-> Async Foreign Scan on pgbench_accounts_2_fdw pgbench_accounts_3  
-> Seq Scan on pgbench_accounts_3 pgbench_accounts_4  
-> Async Foreign Scan on pgbench_accounts_4_fdw pgbench_accounts_5  
-> Async Foreign Scan on pgbench_accounts_5_fdw pgbench_accounts_6  
-> Seq Scan on pgbench_accounts_6 pgbench_accounts_7  
-> Async Foreign Scan on pgbench_accounts_7_fdw pgbench_accounts_8  
-> Async Foreign Scan on pgbench_accounts_8_fdw pgbench_accounts_9  
-> Seq Scan on pgbench_accounts_9 pgbench_accounts_10  
-> Async Foreign Scan on pgbench_accounts_10_fdw pgbench_accounts_11  
-> Async Foreign Scan on pgbench_accounts_11_fdw pgbench_accounts_12  
-> Seq Scan on pgbench_accounts_12 pgbench_accounts_13  
-> Async Foreign Scan on pgbench_accounts_13_fdw pgbench_accounts_14  
-> Async Foreign Scan on pgbench_accounts_14_fdw pgbench_accounts_15  
-> Seq Scan on pgbench_accounts_15 pgbench_accounts_16  
-> Async Foreign Scan on pgbench_accounts_16_fdw pgbench_accounts_17  
-> Async Foreign Scan on pgbench_accounts_17_fdw pgbench_accounts_18  
-> Seq Scan on pgbench_accounts_18 pgbench_accounts_19  
-> Async Foreign Scan on pgbench_accounts_19_fdw pgbench_accounts_20
```

Here we see a typical asynchronous plan. There are asynchronous foreign scans and local sequential scans, which are executed synchronously.

53.2.3. Fetch-all Fallback

There are a lot of cases when operations on data cannot be executed remotely (for example, when some non-immutable function is used in filters, when several sharded tables are joined by an attribute that is not a sharding key, when pushdown of a particular join type is not supported) or when the planner considers local execution to be cheaper. In such cases different operations (selection, joins or aggregations) are not pushed down, but executed locally. This can lead to inefficient query execution due to large inter-cluster traffic and high processing cost on a coordinator. When this happens, you should check if an optimizer has fresh statistics, consider rewriting the query to benefit from different forms of pushdown or at least check that the suggested query plan is reasonable enough.

53.3. How Connections Are Established

Postgres Pro implements a “process per user” client/server model. In this model, every *client process* connects to exactly one *backend process*. As we do not know ahead of time how many connections will be made, we have to use a “supervisor process” that spawns a new backend process every time a connection is requested. This supervisor process is called *postmaster* and listens at a specified TCP/IP port for incoming connections. Whenever it detects a request for a connection, it spawns a new backend process. Those backend processes communicate with each other and with other processes of the *instance* using *semaphores* and *shared memory* to ensure data integrity throughout concurrent data access.

The client process can be any program that understands the Postgres Pro protocol described in [Chapter 56](#). Many clients are based on the C-language library libpq, but several independent implementations of the protocol exist, such as the Java JDBC driver.

Once a connection is established, the client process can send a query to the backend process it's connected to. The query is transmitted using plain text, i.e., there is no parsing done in the client. The backend process parses the query, creates an *execution plan*, executes the plan, and returns the retrieved rows to the client by transmitting them over the established connection.

53.4. The Parser Stage

The *parser stage* consists of two parts:

- The *parser* defined in `gram.y` and `scan.l` is built using the Unix tools `bison` and `flex`.
- The *transformation process* does modifications and augmentations to the data structures returned by the parser.

53.4.1. Parser

The parser has to check the query string (which arrives as plain text) for valid syntax. If the syntax is correct a *parse tree* is built up and handed back; otherwise an error is returned. The parser and lexer are implemented using the well-known Unix tools `bison` and `flex`.

The *lexer* is defined in the file `scan.l` and is responsible for recognizing *identifiers*, the *SQL key words* etc. For every key word or identifier that is found, a *token* is generated and handed to the parser.

The parser is defined in the file `gram.y` and consists of a set of *grammar rules* and *actions* that are executed whenever a rule is fired. The code of the actions (which is actually C code) is used to build up the parse tree.

The file `scan.l` is transformed to the C source file `scan.c` using the program `flex` and `gram.y` is transformed to `gram.c` using `bison`. After these transformations have taken place a normal C compiler can be used to create the parser. Never make any changes to the generated C files as they will be overwritten the next time `flex` or `bison` is called.

Note

The mentioned transformations and compilations are normally done automatically using the *make-files* shipped with the Postgres Pro source distribution.

A detailed description of `bison` or the grammar rules given in `gram.y` would be beyond the scope of this manual. There are many books and documents dealing with `flex` and `bison`. You should be familiar with `bison` before you start to study the grammar given in `gram.y` otherwise you won't understand what happens there.

53.4.2. Transformation Process

The parser stage creates a parse tree using only fixed rules about the syntactic structure of SQL. It does not make any lookups in the system catalogs, so there is no possibility to understand the detailed semantics of the requested operations. After the parser completes, the *transformation process* takes the tree handed back by the parser as input and does the semantic interpretation needed to understand which tables, functions, and operators are referenced by the query. The data structure that is built to represent this information is called the *query tree*.

The reason for separating raw parsing from semantic analysis is that system catalog lookups can only be done within a transaction, and we do not wish to start a transaction immediately upon receiving a query string. The raw parsing stage is sufficient to identify the transaction control commands (`BEGIN`, `ROLLBACK`, etc.), and these can then be correctly executed without any further analysis. Once we know that we are dealing with an actual query (such as `SELECT` or `UPDATE`), it is okay to start a transaction if we're not already in one. Only then can the transformation process be invoked.

The query tree created by the transformation process is structurally similar to the raw parse tree in most places, but it has many differences in detail. For example, a `FuncCall` node in the parse tree represents something that looks syntactically like a function call. This might be transformed to either a `FuncExpr` or `Aggref` node depending on whether the referenced name turns out to be an ordinary function or an

aggregate function. Also, information about the actual data types of columns and expression results is added to the query tree.

53.5. The Postgres Pro Rule System

Postgres Pro supports a powerful *rule system* for the specification of *views* and ambiguous *view updates*. Originally the Postgres Pro rule system consisted of two implementations:

- The first one worked using *row level* processing and was implemented deep in the *executor*. The rule system was called whenever an individual row had been accessed. This implementation was removed in 1995 when the last official release of the Berkeley Postgres project was transformed into Postgres95.
- The second implementation of the rule system is a technique called *query rewriting*. The *rewrite system* is a module that exists between the *parser stage* and the *planner/optimizer*. This technique is still implemented.

The query rewriter is discussed in some detail in [Chapter 42](#), so there is no need to cover it here. We will only point out that both the input and the output of the rewriter are query trees, that is, there is no change in the representation or level of semantic detail in the trees. Rewriting can be thought of as a form of macro expansion.

53.6. Planner/Optimizer

The task of the *planner/optimizer* is to create an optimal execution plan. A given SQL query (and hence, a query tree) can be actually executed in a wide variety of different ways, each of which will produce the same set of results. If it is computationally feasible, the query optimizer will examine each of these possible execution plans, ultimately selecting the execution plan that is expected to run the fastest.

Note

In some situations, examining each possible way in which a query can be executed would take an excessive amount of time and memory. In particular, this occurs when executing queries involving large numbers of join operations. In order to determine a reasonable (not necessarily optimal) query plan in a reasonable amount of time, Postgres Pro uses a *Genetic Query Optimizer* (see [Chapter 61](#)) when the number of joins exceeds a threshold (see [geqo_threshold](#)).

The planner's search procedure actually works with data structures called *paths*, which are simply cut-down representations of plans containing only as much information as the planner needs to make its decisions. After the cheapest path is determined, a full-fledged *plan tree* is built to pass to the executor. This represents the desired execution plan in sufficient detail for the executor to run it. In the rest of this section we'll ignore the distinction between paths and plans.

53.6.1. Generating Possible Plans

The planner/optimizer starts by generating plans for scanning each individual relation (table) used in the query. The possible plans are determined by the available indexes on each relation. There is always the possibility of performing a sequential scan on a relation, so a sequential scan plan is always created. Assume an index is defined on a relation (for example a B-tree index) and a query contains the restriction `relation.attribute OPR constant`. If `relation.attribute` happens to match the key of the B-tree index and `OPR` is one of the operators listed in the index's *operator class*, another plan is created using the B-tree index to scan the relation. If there are further indexes present and the restrictions in the query happen to match a key of an index, further plans will be considered. Index scan plans are also generated for indexes that have a sort ordering that can match the query's `ORDER BY` clause (if any), or a sort ordering that might be useful for merge joining (see below).

If the query requires joining two or more relations, plans for joining relations are considered after all feasible plans have been found for scanning single relations. The three available join strategies are:

- *nested loop join*: The right relation is scanned once for every row found in the left relation. This strategy is easy to implement but can be very time consuming. (However, if the right relation can be scanned with an index scan, this can be a good strategy. It is possible to use values from the current row of the left relation as keys for the index scan of the right.)
- *merge join*: Each relation is sorted on the join attributes before the join starts. Then the two relations are scanned in parallel, and matching rows are combined to form join rows. This kind of join is attractive because each relation has to be scanned only once. The required sorting might be achieved either by an explicit sort step, or by scanning the relation in the proper order using an index on the join key.
- *hash join*: the right relation is first scanned and loaded into a hash table, using its join attributes as hash keys. Next the left relation is scanned and the appropriate values of every row found are used as hash keys to locate the matching rows in the table.

When the query involves more than two relations, the final result must be built up by a tree of join steps, each with two inputs. The planner examines different possible join sequences to find the cheapest one.

If the query uses fewer than `geqo_threshold` relations, a near-exhaustive search is conducted to find the best join sequence. The planner preferentially considers joins between any two relations for which there exists a corresponding join clause in the `WHERE` qualification (i.e., for which a restriction like `where rel1.attr1=rel2.attr2` exists). Join pairs with no join clause are considered only when there is no other choice, that is, a particular relation has no available join clauses to any other relation. All possible plans are generated for every join pair considered by the planner, and the one that is (estimated to be) the cheapest is chosen.

When `geqo_threshold` is exceeded, the join sequences considered are determined by heuristics, as described in [Chapter 61](#). Otherwise the process is the same.

The finished plan tree consists of sequential or index scans of the base relations, plus nested-loop, merge, or hash join nodes as needed, plus any auxiliary steps needed, such as sort nodes or aggregate-function calculation nodes. Most of these plan node types have the additional ability to do *selection* (discarding rows that do not meet a specified Boolean condition) and *projection* (computation of a derived column set based on given column values, that is, evaluation of scalar expressions where needed). One of the responsibilities of the planner is to attach selection conditions from the `WHERE` clause and computation of required output expressions to the most appropriate nodes of the plan tree.

53.7. Executor

The *executor* takes the plan created by the planner/optimizer and recursively processes it to extract the required set of rows. This is essentially a demand-pull pipeline mechanism. Each time a plan node is called, it must deliver one more row, or report that it is done delivering rows.

To provide a concrete example, assume that the top node is a `MergeJoin` node. Before any merge can be done two rows have to be fetched (one from each subplan). So the executor recursively calls itself to process the subplans (it starts with the subplan attached to `lefttree`). The new top node (the top node of the left subplan) is, let's say, a `Sort` node and again recursion is needed to obtain an input row. The child node of the `Sort` might be a `SeqScan` node, representing actual reading of a table. Execution of this node causes the executor to fetch a row from the table and return it up to the calling node. The `Sort` node will repeatedly call its child to obtain all the rows to be sorted. When the input is exhausted (as indicated by the child node returning a `NULL` instead of a row), the `Sort` code performs the sort, and finally is able to return its first output row, namely the first one in sorted order. It keeps the remaining rows stored so that it can deliver them in sorted order in response to later demands.

The `MergeJoin` node similarly demands the first row from its right subplan. Then it compares the two rows to see if they can be joined; if so, it returns a join row to its caller. On the next call, or immediately if it cannot join the current pair of inputs, it advances to the next row of one table or the other (depending on how the comparison came out), and again checks for a match. Eventually, one subplan or the other is exhausted, and the `MergeJoin` node returns `NULL` to indicate that no more join rows can be formed.

Complex queries can involve many levels of plan nodes, but the general approach is the same: each node computes and returns its next output row each time it is called. Each node is also responsible for applying any selection or projection expressions that were assigned to it by the planner.

The executor mechanism is used to evaluate all five basic SQL query types: `SELECT`, `INSERT`, `UPDATE`, `DELETE`, and `MERGE`. For `SELECT`, the top-level executor code only needs to send each row returned by the query plan tree off to the client. `INSERT ... SELECT`, `UPDATE`, `DELETE`, and `MERGE` are effectively `SELECT`s under a special top-level plan node called `ModifyTable`.

`INSERT ... SELECT` feeds the rows up to `ModifyTable` for insertion. For `UPDATE`, the planner arranges that each computed row includes all the updated column values, plus the *TID* (tuple ID, or row ID) of the original target row; this data is fed up to the `ModifyTable` node, which uses the information to create a new updated row and mark the old row deleted. For `DELETE`, the only column that is actually returned by the plan is the TID, and the `ModifyTable` node simply uses the TID to visit each target row and mark it deleted. For `MERGE`, the planner joins the source and target relations, and includes all column values required by any of the `WHEN` clauses, plus the TID of the target row; this data is fed up to the `ModifyTable` node, which uses the information to work out which `WHEN` clause to execute, and then inserts, updates or deletes the target row, as required.

A simple `INSERT ... VALUES` command creates a trivial plan tree consisting of a single `Result` node, which computes just one result row, feeding that up to `ModifyTable` to perform the insertion.

Chapter 54. System Catalogs

The system catalogs are the place where a relational database management system stores schema meta-data, such as information about tables and columns, and internal bookkeeping information. Postgres Pro's system catalogs are regular tables. You can drop and recreate the tables, add columns, insert and update values, and severely mess up your system that way. Normally, one should not change the system catalogs by hand, there are normally SQL commands to do that. (For example, `CREATE DATABASE` inserts a row into the `pg_database` catalog — and actually creates the database on disk.) There are some exceptions for particularly esoteric operations, but many of those have been made available as SQL commands over time, and so the need for direct manipulation of the system catalogs is ever decreasing.

54.1. Overview

Table 54.1 lists the system catalogs. More detailed documentation of each catalog follows below.

Most system catalogs are copied from the template database during database creation and are thereafter database-specific. A few catalogs are physically shared across all databases in a cluster; these are noted in the descriptions of the individual catalogs.

Table 54.1. System Catalogs

Catalog Name	Purpose
<code>pg_aggregate</code>	aggregate functions
<code>pg_am</code>	relation access methods
<code>pg_amop</code>	access method operators
<code>pg_amproc</code>	access method support functions
<code>pg_attrdef</code>	column default values
<code>pg_attribute</code>	table columns (“attributes”)
<code>pg_authid</code>	authorization identifiers (roles)
<code>pg_auth_members</code>	authorization identifier membership relationships
<code>pg_cast</code>	casts (data type conversions)
<code>pg_class</code>	tables, indexes, sequences, views (“relations”)
<code>pg_collation</code>	collations (locale information)
<code>pg_constraint</code>	check constraints, unique constraints, primary key constraints, foreign key constraints
<code>pg_conversion</code>	encoding conversion information
<code>pg_database</code>	databases within this database cluster
<code>pg_db_role_setting</code>	per-role and per-database settings
<code>pg_default_acl</code>	default privileges for object types
<code>pg_depend</code>	dependencies between database objects
<code>pg_description</code>	descriptions or comments on database objects
<code>pg_enum</code>	enum label and value definitions
<code>pg_event_trigger</code>	event triggers
<code>pg_extension</code>	installed extensions
<code>pg_foreign_data_wrapper</code>	foreign-data wrapper definitions
<code>pg_foreign_server</code>	foreign server definitions
<code>pg_foreign_table</code>	additional foreign table information
<code>pg_index</code>	additional index information

Catalog Name	Purpose
pg_inherits	table inheritance hierarchy
pg_init_privs	object initial privileges
pg_language	languages for writing functions
pg_largeobject	data pages for large objects
pg_largeobject_metadata	metadata for large objects
pg_namespace	schemas
pg_opclass	access method operator classes
pg_operator	operators
pg_opfamily	access method operator families
pg_parameter_acl	configuration parameters for which privileges have been granted
pg_partitioned_table	information about partition key of tables
pg_policy	row-security policies
pg_proc	functions and procedures
pg_profile	profiles, a set of authentication restrictions
pg_publication	publications for logical replication
pg_publication_namespace	schema to publication mapping
pg_publication_rel	relation to publication mapping
pg_range	information about range types
pg_replication_origin	registered replication origins
pg_rewrite	query rewrite rules
pg_role_password	roles password history
pg_seclabel	security labels on database objects
pg_sequence	information about sequences
pg_shdepend	dependencies on shared objects
pg_shdescription	comments on shared objects
pg_shseclabel	security labels on shared database objects
pg_statistic	planner statistics
pg_statistic_ext	extended planner statistics (definition)
pg_statistic_ext_data	extended planner statistics (built statistics)
pg_subscription	logical replication subscriptions
pg_subscription_rel	relation state for subscriptions
pg_tablespace	tablespaces within this database cluster
pg_transform	transforms (data type to procedural language conversions)
pg_trigger	triggers
pg_ts_config	text search configurations
pg_ts_config_map	text search configurations' token mappings
pg_ts_dict	text search dictionaries
pg_ts_parser	text search parsers

Catalog Name	Purpose
pg_ts_template	text search templates
pg_type	data types
pg_user_mapping	mappings of users to foreign servers

54.2. pg_aggregate

The catalog `pg_aggregate` stores information about aggregate functions. An aggregate function is a function that operates on a set of values (typically one column from each row that matches a query condition) and returns a single value computed from all these values. Typical aggregate functions are `sum`, `count`, and `max`. Each entry in `pg_aggregate` is an extension of an entry in [pg_proc](#). The `pg_proc` entry carries the aggregate's name, input and output data types, and other information that is similar to ordinary functions.

Table 54.2. `pg_aggregate` Columns

Column Type	Description
<code>aggfnoid regproc</code> (references pg_proc .oid)	<code>pg_proc</code> OID of the aggregate function
<code>aggkind char</code>	Aggregate kind: <code>n</code> for “normal” aggregates, <code>o</code> for “ordered-set” aggregates, or <code>h</code> for “hypothetical-set” aggregates
<code>aggnumdirectargs int2</code>	Number of direct (non-aggregated) arguments of an ordered-set or hypothetical-set aggregate, counting a variadic array as one argument. If equal to <code>pronargs</code> , the aggregate must be variadic and the variadic array describes the aggregated arguments as well as the final direct arguments. Always zero for normal aggregates.
<code>aggtransfn regproc</code> (references pg_proc .oid)	Transition function
<code>aggfinalfn regproc</code> (references pg_proc .oid)	Final function (zero if none)
<code>aggcombinefn regproc</code> (references pg_proc .oid)	Combine function (zero if none)
<code>aggserialfn regproc</code> (references pg_proc .oid)	Serialization function (zero if none)
<code>aggdeserialfn regproc</code> (references pg_proc .oid)	Deserialization function (zero if none)
<code>aggmtransfn regproc</code> (references pg_proc .oid)	Forward transition function for moving-aggregate mode (zero if none)
<code>aggminvtransfn regproc</code> (references pg_proc .oid)	Inverse transition function for moving-aggregate mode (zero if none)
<code>aggmfinalfn regproc</code> (references pg_proc .oid)	Final function for moving-aggregate mode (zero if none)
<code>aggfinalextra bool</code>	True to pass extra dummy arguments to <code>aggfinalfn</code>
<code>aggmfinalextra bool</code>	True to pass extra dummy arguments to <code>aggmfinalfn</code>
<code>aggfinalmodify char</code>	Whether <code>aggfinalfn</code> modifies the transition state value: <code>r</code> if it is read-only, <code>s</code> if the <code>aggtransfn</code> cannot be applied after the <code>aggfinalfn</code> , or <code>w</code> if it writes on the value

Column Type	Description
aggmfinalmodify char	Like aggfinalmodify, but for the aggmfinalfn
aggstortop oid (references pg_operator .oid)	Associated sort operator (zero if none)
aggtranstype oid (references pg_type .oid)	Data type of the aggregate function's internal transition (state) data
aggtransspace int4	Approximate average size (in bytes) of the transition state data, or zero to use a default estimate
aggmtranstype oid (references pg_type .oid)	Data type of the aggregate function's internal transition (state) data for moving-aggregate mode (zero if none)
aggmtransspace int4	Approximate average size (in bytes) of the transition state data for moving-aggregate mode, or zero to use a default estimate
agginitval text	The initial value of the transition state. This is a text field containing the initial value in its external string representation. If this field is null, the transition state value starts out null.
aggminitval text	The initial value of the transition state for moving-aggregate mode. This is a text field containing the initial value in its external string representation. If this field is null, the transition state value starts out null.

New aggregate functions are registered with the [CREATE AGGREGATE](#) command. See [Section 39.12](#) for more information about writing aggregate functions and the meaning of the transition functions, etc.

54.3. pg_am

The catalog `pg_am` stores information about relation access methods. There is one row for each access method supported by the system. Currently, only tables and indexes have access methods. The requirements for table and index access methods are discussed in detail in [Chapter 62](#) and [Chapter 63](#) respectively.

Table 54.3. `pg_am` Columns

Column Type	Description
oid oid	Row identifier
amname name	Name of the access method
amhandler regproc (references pg_proc .oid)	OID of a handler function that is responsible for supplying information about the access method
amtype char	t = table (including materialized views), i = index.

Note

Before Postgres Pro 9.6, `pg_am` contained many additional columns representing properties of index access methods. That data is now only directly visible at the C code level. However, `pg_in-`

`dex_column_has_property()` and related functions have been added to allow SQL queries to inspect index access method properties; see [Table 9.74](#).

54.4. `pg_amop`

The catalog `pg_amop` stores information about operators associated with access method operator families. There is one row for each operator that is a member of an operator family. A family member can be either a *search* operator or an *ordering* operator. An operator can appear in more than one family, but cannot appear in more than one search position nor more than one ordering position within a family. (It is allowed, though unlikely, for an operator to be used for both search and ordering purposes.)

Table 54.4. `pg_amop` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>amopfamily oid (references <code>pg_opfamily</code> .oid)</code>	The operator family this entry is for
<code>amoplefttype oid (references <code>pg_type</code> .oid)</code>	Left-hand input data type of operator
<code>amoprightrighttype oid (references <code>pg_type</code> .oid)</code>	Right-hand input data type of operator
<code>amopstrategy int2</code>	Operator strategy number
<code>amoppurpose char</code>	Operator purpose, either <code>s</code> for search or <code>o</code> for ordering
<code>amopopr oid (references <code>pg_operator</code> .oid)</code>	OID of the operator
<code>amopmethod oid (references <code>pg_am</code> .oid)</code>	Index access method operator family is for
<code>amopsortfamily oid (references <code>pg_opfamily</code> .oid)</code>	The B-tree operator family this entry sorts according to, if an ordering operator; zero if a search operator

A “search” operator entry indicates that an index of this operator family can be searched to find all rows satisfying `WHERE indexed_column operator constant`. Obviously, such an operator must return `boolean`, and its left-hand input type must match the index's column data type.

An “ordering” operator entry indicates that an index of this operator family can be scanned to return rows in the order represented by `ORDER BY indexed_column operator constant`. Such an operator could return any sortable data type, though again its left-hand input type must match the index's column data type. The exact semantics of the `ORDER BY` are specified by the `amopsortfamily` column, which must reference a B-tree operator family for the operator's result type.

Note

At present, it's assumed that the sort order for an ordering operator is the default for the referenced operator family, i.e., `ASC NULLS LAST`. This might someday be relaxed by adding additional columns to specify sort options explicitly.

An entry's `amopmethod` must match the `opfmeth` of its containing operator family (including `amopmethod` here is an intentional denormalization of the catalog structure for performance reasons). Al-

so, `amoplefttype` and `amoprightright` must match the `oprleft` and `oprright` fields of the referenced `pg_operator` entry.

54.5. `pg_amproc`

The catalog `pg_amproc` stores information about support functions associated with access method operator families. There is one row for each support function belonging to an operator family.

Table 54.5. `pg_amproc` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>amprocfamily oid</code> (references <code>pg_opfamily .oid</code>)	The operator family this entry is for
<code>amproclefttype oid</code> (references <code>pg_type .oid</code>)	Left-hand input data type of associated operator
<code>amprocrighttype oid</code> (references <code>pg_type .oid</code>)	Right-hand input data type of associated operator
<code>amprocnum int2</code>	Support function number
<code>amproc regproc</code> (references <code>pg_proc .oid</code>)	OID of the function

The usual interpretation of the `amproclefttype` and `amprocrighttype` fields is that they identify the left and right input types of the operator(s) that a particular support function supports. For some access methods these match the input data type(s) of the support function itself, for others not. There is a notion of “default” support functions for an index, which are those with `amproclefttype` and `amprocrighttype` both equal to the index operator class's `opcintype`.

54.6. `pg_attrdef`

The catalog `pg_attrdef` stores column default values. The main information about columns is stored in `pg_attribute`. Only columns for which a default value has been explicitly set will have an entry here.

Table 54.6. `pg_attrdef` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>adrelid oid</code> (references <code>pg_class .oid</code>)	The table this column belongs to
<code>adnum int2</code> (references <code>pg_attribute .attnum</code>)	The number of the column
<code>adbin pg_node_tree</code>	The column default value, in <code>nodeToString()</code> representation. Use <code>pg_get_expr(adbin, adrelid)</code> to convert it to an SQL expression.

54.7. `pg_attribute`

The catalog `pg_attribute` stores information about table columns. There will be exactly one `pg_attribute` row for every column in every table in the database. (There will also be attribute entries for indexes, and indeed all objects that have `pg_class` entries.)

The term attribute is equivalent to column and is used for historical reasons.

Table 54.7. pg_attribute Columns

Column Type	Description
attrelid oid (references <code>pg_class</code> .oid)	The table this column belongs to
attname name	The column name
atttypid oid (references <code>pg_type</code> .oid)	The data type of this column (zero for a dropped column)
attlen int2	A copy of <code>pg_type.typelen</code> of this column's type
attnum int2	The number of the column. Ordinary columns are numbered from 1 up. System columns, such as <code>ctid</code> , have (arbitrary) negative numbers.
attcacheoff int4	Always -1 in storage, but when loaded into a row descriptor in memory this might be updated to cache the offset of the attribute within the row
atttypmod int4	<code>atttypmod</code> records type-specific data supplied at table creation time (for example, the maximum length of a <code>varchar</code> column). It is passed to type-specific input functions and length coercion functions. The value will generally be -1 for types that do not need <code>atttypmod</code> .
attnndims int2	Number of dimensions, if the column is an array type; otherwise 0. (Presently, the number of dimensions of an array is not enforced, so any nonzero value effectively means "it's an array".)
attbyval bool	A copy of <code>pg_type.typbyval</code> of this column's type
attalign char	A copy of <code>pg_type.typalign</code> of this column's type
attstorage char	Normally a copy of <code>pg_type.typstorage</code> of this column's type. For TOAST-able data types, this can be altered after column creation to control storage policy.
attcompression char	The current compression method of the column. Typically this is <code>'\0'</code> to specify use of the current default setting (see default_toast_compression). Otherwise, <code>'p'</code> selects <code>pglz</code> compression, while <code>'l'</code> selects <code>LZ4</code> compression. However, this field is ignored whenever <code>attstorage</code> does not allow compression.
attnotnull bool	This represents a not-null constraint.
atthasdef bool	This column has a default expression or generation expression, in which case there will be a corresponding entry in the <code>pg_attrdef</code> catalog that actually defines the expression. (Check <code>attgenerated</code> to determine whether this is a default or a generation expression.)
atthasmissing bool	This column has a value which is used where the column is entirely missing from the row, as happens when a column is added with a non-volatile <code>DEFAULT</code> value after the row is created. The actual value used is stored in the <code>attmissingval</code> column.
attidentity char	

Column Type	Description
	If a zero byte (' '), then not an identity column. Otherwise, <code>a</code> = generated always, <code>d</code> = generated by default.
<code>attgenerated char</code>	If a zero byte (' '), then not a generated column. Otherwise, <code>s</code> = stored. (Other values might be added in the future.)
<code>attisdropped bool</code>	This column has been dropped and is no longer valid. A dropped column is still physically present in the table, but is ignored by the parser and so cannot be accessed via SQL.
<code>attislocal bool</code>	This column is defined locally in the relation. Note that a column can be locally defined and inherited simultaneously.
<code>attinhcount int2</code>	The number of direct ancestors this column has. A column with a nonzero number of ancestors cannot be dropped nor renamed.
<code>attcollation oid</code> (references <code>pg_collation .oid</code>)	The defined collation of the column, or zero if the column is not of a collatable data type
<code>attstattarget int2</code>	<code>attstattarget</code> controls the level of detail of statistics accumulated for this column by ANALYZE . A zero value indicates that no statistics should be collected. A null value says to use the system default statistics target. The exact meaning of positive values is data type-dependent. For scalar data types, <code>attstattarget</code> is both the target number of “most common values” to collect, and the target number of histogram bins to create.
<code>attacl aclitem[]</code>	Column-level access privileges, if any have been granted specifically on this column
<code>attoptions text[]</code>	Attribute-level options, as “keyword=value” strings
<code>attfdwoptions text[]</code>	Attribute-level foreign data wrapper options, as “keyword=value” strings
<code>attmissingval anyarray</code>	This column has a one element array containing the value used when the column is entirely missing from the row, as happens when the column is added with a non-volatile <code>DEFAULT</code> value after the row is created. The value is only used when <code>atthasmissing</code> is true. If there is no value the column is null.

In a dropped column's `pg_attribute` entry, `atttypid` is reset to zero, but `attlen` and the other fields copied from `pg_type` are still valid. This arrangement is needed to cope with the situation where the dropped column's data type was later dropped, and so there is no `pg_type` row anymore. `attlen` and the other fields can be used to interpret the contents of a row of the table.

54.8. pg_authid

The catalog `pg_authid` contains information about database authorization identifiers (roles). A role subsumes the concepts of “users” and “groups”. A user is essentially just a role with the `rolcanlogin` flag set. Any role (with or without `rolcanlogin`) can have other roles as members; see [pg_auth_members](#).

Since this catalog contains passwords, it must not be publicly readable. `pg_roles` is a publicly readable view on `pg_authid` that blanks out the password field.

[Chapter 21](#) contains detailed information about user and privilege management.

Because user identities are cluster-wide, `pg_authid` is shared across all databases of a cluster: there is only one copy of `pg_authid` per cluster, not one per database.

Table 54.8. pg_authid Columns

Column	Type	Description
oid	oid	Row identifier
rolname	name	Role name
rolsuper	bool	Role has superuser privileges
rolinherit	bool	Role automatically inherits privileges of roles it is a member of
rolcreatorole	bool	Role can create more roles
rolcreatedb	bool	Role can create databases
rolcanlogin	bool	Role can log in. That is, this role can be given as the initial session authorization identifier.
rolreplication	bool	Role is a replication role. A replication role can initiate replication connections and create and drop replication slots.
rolbypassrsls	bool	Role bypasses every row-level security policy, see Section 5.9 for more information.
rolconnlimit	int4	For roles that can log in, this sets maximum number of concurrent connections this role can make. -1 means no limit.
rolprofile	regprofile (references pg_profile .oid)	The OID of the role's profile
rolloginattempts	int4	Number of consecutive failed login attempts of a user. It is always 0 if FAILED_LOGIN_ATTEMPTS value of the corresponding role profile is UNLIMITED (see CREATE PROFILE). After a successful login, rolloginattempts is reset to 0.
rollastlogin	timestampz	The timestamp the role logged in last time
rolfirstfailedauth	timestampz	The timestamp of the role's first authentication failure
rolstatus	int2	Status of the role: 0 for the active role, 1 the role is manually locked (see ACCOUNT LOCK in ALTER ROLE), 2 the role is locked because of inactivity (see parameter USER_INACTIVE_TIME in CREATE PROFILE), 4 the role is locked because the number of consecutive authentication failures has reached the limit (see parameter FAILED_LOGIN_ATTEMPTS in CREATE PROFILE)
rolpassword	text	Password (possibly encrypted); null if none. The format depends on the form of encryption used.
rolvaliduntil	timestampz	Password expiry time (only used for password authentication); null if no expiration
rolpasssetat	timestampz	Password set time (only used for password authentication); null if password is not set.

For an MD5 encrypted password, `rolpassword` column will begin with the string `md5` followed by a 32-character hexadecimal MD5 hash. The MD5 hash will be of the user's password concatenated to their user name. For example, if user `joe` has password `xyzyzy`, Postgres Pro will store the md5 hash of `xyzyzyjoe`.

If the password is encrypted with SCRAM-SHA-256, it has the format:

```
SCRAM-SHA-256$<iteration count>:<salt>$<StoredKey>:<ServerKey>
```

where `salt`, `StoredKey` and `ServerKey` are in Base64 encoded format. This format is the same as that specified by [RFC 5803](#).

A password that does not follow either of those formats is assumed to be unencrypted.

54.9. pg_auth_members

The catalog `pg_auth_members` shows the membership relations between roles. Any non-circular set of relationships is allowed.

Because user identities are cluster-wide, `pg_auth_members` is shared across all databases of a cluster: there is only one copy of `pg_auth_members` per cluster, not one per database.

Table 54.9. pg_auth_members Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>roleid oid (references pg_authid .oid)</code>	ID of a role that has a member
<code>member oid (references pg_authid .oid)</code>	ID of a role that is a member of <code>roleid</code>
<code>grantor oid (references pg_authid .oid)</code>	ID of the role that granted this membership
<code>admin_option bool</code>	True if member can grant membership in <code>roleid</code> to others
<code>inherit_option bool</code>	True if the member automatically inherits the privileges of the granted role
<code>set_option bool</code>	True if the member can <code>SET ROLE</code> to the granted role

54.10. pg_cast

The catalog `pg_cast` stores data type conversion paths, both built-in and user-defined.

It should be noted that `pg_cast` does not represent every type conversion that the system knows how to perform; only those that cannot be deduced from some generic rule. For example, casting between a domain and its base type is not explicitly represented in `pg_cast`. Another important exception is that “automatic I/O conversion casts”, those performed using a data type's own I/O functions to convert to or from `text` or other string types, are not explicitly represented in `pg_cast`.

Table 54.10. pg_cast Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>castsource oid (references pg_type .oid)</code>	

Column Type	Description
	OID of the source data type
<code>casttarget oid</code> (references <code>pg_type .oid</code>)	OID of the target data type
<code>castfunc oid</code> (references <code>pg_proc .oid</code>)	The OID of the function to use to perform this cast. Zero is stored if the cast method doesn't require a function.
<code>castcontext char</code>	Indicates what contexts the cast can be invoked in. <code>e</code> means only as an explicit cast (using <code>CAST</code> or <code>::</code> syntax). <code>a</code> means implicitly in assignment to a target column, as well as explicitly. <code>i</code> means implicitly in expressions, as well as the other cases.
<code>castmethod char</code>	Indicates how the cast is performed. <code>f</code> means that the function specified in the <code>castfunc</code> field is used. <code>i</code> means that the input/output functions are used. <code>b</code> means that the types are binary-coercible, thus no conversion is required.

The cast functions listed in `pg_cast` must always take the cast source type as their first argument type, and return the cast destination type as their result type. A cast function can have up to three arguments. The second argument, if present, must be type `integer`; it receives the type modifier associated with the destination type, or -1 if there is none. The third argument, if present, must be type `boolean`; it receives `true` if the cast is an explicit cast, `false` otherwise.

It is legitimate to create a `pg_cast` entry in which the source and target types are the same, if the associated function takes more than one argument. Such entries represent “length coercion functions” that coerce values of the type to be legal for a particular type modifier value.

When a `pg_cast` entry has different source and target types and a function that takes more than one argument, it represents converting from one type to another and applying a length coercion in a single step. When no such entry is available, coercion to a type that uses a type modifier involves two steps, one to convert between data types and a second to apply the modifier.

54.11. `pg_class`

The catalog `pg_class` describes tables and other objects that have columns or are otherwise similar to a table. This includes indexes (but see also `pg_index`), sequences (but see also `pg_sequence`), views, materialized views, composite types, and TOAST tables; see `relkind`. Below, when we mean all of these kinds of objects we speak of “relations”. Not all of `pg_class`'s columns are meaningful for all relation kinds.

Table 54.11. `pg_class` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>relname name</code>	Name of the table, index, view, etc.
<code>relnamespace oid</code> (references <code>pg_namespace .oid</code>)	The OID of the namespace that contains this relation
<code>reltype oid</code> (references <code>pg_type .oid</code>)	The OID of the data type that corresponds to this table's row type, if any; zero for indexes, sequences, and toast tables, which have no <code>pg_type</code> entry
<code>reloftype oid</code> (references <code>pg_type .oid</code>)	For typed tables, the OID of the underlying composite type; zero for all other relations

Column Type	Description
<code>relowner oid</code> (references <code>pg_authid .oid</code>)	Owner of the relation
<code>relam oid</code> (references <code>pg_am .oid</code>)	The access method used to access this table or index. Not meaningful if the relation is a sequence or has no on-disk file, except for partitioned tables, where, if set, it takes precedence over <code>default_table_access_method</code> when determining the access method to use for partitions created when one is not specified in the creation command.
<code>relfilenode oid</code>	Name of the on-disk file of this relation; zero means this is a “mapped” relation whose disk file name is determined by low-level state
<code>reltablespace oid</code> (references <code>pg_tablespace .oid</code>)	The tablespace in which this relation is stored. If zero, the database's default tablespace is implied. Not meaningful if the relation has no on-disk file, except for partitioned tables, where this is the tablespace in which partitions will be created when one is not specified in the creation command.
<code>relpages int4</code>	Size of the on-disk representation of this table in pages (of size <code>BLCKSZ</code>). This is only an estimate used by the planner. It is updated by <code>VACUUM</code> , <code>ANALYZE</code> , and a few DDL commands such as <code>CREATE INDEX</code> .
<code>reltuples float4</code>	Number of live rows in the table. This is only an estimate used by the planner. It is updated by <code>VACUUM</code> , <code>ANALYZE</code> , and a few DDL commands such as <code>CREATE INDEX</code> . If the table has never yet been vacuumed or analyzed, <code>reltuples</code> contains <code>-1</code> indicating that the row count is unknown.
<code>relallvisible int4</code>	Number of pages that are marked all-visible in the table's visibility map. This is only an estimate used by the planner. It is updated by <code>VACUUM</code> , <code>ANALYZE</code> , and a few DDL commands such as <code>CREATE INDEX</code> .
<code>reltoastrelid oid</code> (references <code>pg_class .oid</code>)	OID of the TOAST table associated with this table, zero if none. The TOAST table stores large attributes “out of line” in a secondary table.
<code>relhasindex bool</code>	True if this is a table and it has (or recently had) any indexes
<code>relisshared bool</code>	True if this table is shared across all databases in the cluster. Only certain system catalogs (such as <code>pg_database</code>) are shared.
<code>relpersistence char</code> p = permanent table/sequence, u = unlogged table/sequence, t = temporary table/sequence, c = constant table	
<code>relkind char</code> r = ordinary table, i = index, S = sequence, t = TOAST table, v = view, m = materialized view, c = composite type, f = foreign table, p = partitioned table, I = partitioned index	
<code>relnatts int2</code>	Number of user columns in the relation (system columns not counted). There must be this many corresponding entries in <code>pg_attribute</code> . See also <code>pg_attribute .attnum</code> .
<code>relchecks int2</code>	Number of <code>CHECK</code> constraints on the table; see <code>pg_constraint</code> catalog
<code>relhasrules bool</code>	True if table has (or once had) rules; see <code>pg_rewrite</code> catalog

Column Type	Description
<code>relhastriggers</code> bool	True if table has (or once had) triggers; see <code>pg_trigger</code> catalog
<code>relhassubclass</code> bool	True if table or index has (or once had) any inheritance children or partitions
<code>relrowsecurity</code> bool	True if table has row-level security enabled; see <code>pg_policy</code> catalog
<code>relforcerowsecurity</code> bool	True if row-level security (when enabled) will also apply to table owner; see <code>pg_policy</code> catalog
<code>relispopulated</code> bool	True if relation is populated (this is true for all relations other than some materialized views)
<code>relreplident</code> char	Columns used to form “replica identity” for rows: <code>d</code> = default (primary key, if any), <code>n</code> = nothing, <code>f</code> = all columns, <code>i</code> = index with <code>indisreplident</code> set (same as nothing if the index used has been dropped)
<code>relispartition</code> bool	True if table or index is a partition
<code>relrewrite</code> oid (references <code>pg_class</code> .oid)	For new relations being written during a DDL operation that requires a table rewrite, this contains the OID of the original relation; otherwise zero. That state is only visible internally; this field should never contain anything other than zero for a user-visible relation.
<code>relfrozenxid</code> xid	All transaction IDs before this one have been replaced with a permanent (“frozen”) transaction ID in this table. This is used to track whether the table needs to be vacuumed in order to allow <code>pg_xact</code> to be shrunk. Zero (<code>InvalidTransactionId</code>) if the relation is not a table.
<code>relminmxid</code> xid	All multixact IDs before this one have been replaced by a transaction ID in this table. This is used to track whether the table needs to be vacuumed in order to allow <code>pg_multixact</code> to be shrunk. Zero (<code>InvalidMultiXactId</code>) if the relation is not a table.
<code>relacl</code> aclitem[]	Access privileges; see Section 5.8 for details
<code>reloptions</code> text[]	Access-method-specific options, as “keyword=value” strings
<code>relpartbound</code> <code>pg_node_tree</code>	If table is a partition (see <code>relispartition</code>), internal representation of the partition bound

Several of the Boolean flags in `pg_class` are maintained lazily: they are guaranteed to be true if that's the correct state, but may not be reset to false immediately when the condition is no longer true. For example, `relhasindex` is set by `CREATE INDEX`, but it is never cleared by `DROP INDEX`. Instead, `VACUUM` clears `relhasindex` if it finds the table has no indexes. This arrangement avoids race conditions and improves concurrency.

54.12. `pg_collation`

The catalog `pg_collation` describes the available collations, which are essentially mappings from an SQL name to operating system locale categories. See [Section 23.2](#) for more information.

Table 54.12. `pg_collation` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>collname name</code>	Collation name (unique per namespace and encoding)
<code>collnamespace oid</code> (references <code>pg_namespace .oid</code>)	The OID of the namespace that contains this collation
<code>collowner oid</code> (references <code>pg_authid .oid</code>)	Owner of the collation
<code>collprovider char</code>	Provider of the collation: <code>d</code> = database default, <code>b</code> = builtin, <code>c</code> = libc, <code>i</code> = icu
<code>collisdeterministic bool</code>	Is the collation deterministic?
<code>collencoding int4</code>	Encoding in which the collation is applicable, or -1 if it works for any encoding
<code>collcollate text</code>	<code>LC_COLLATE</code> for this collation object. If the provider is not <code>libc</code> , <code>collcollate</code> is NULL and <code>colllocale</code> is used instead.
<code>collctype text</code>	<code>LC_CTYPE</code> for this collation object. If the provider is not <code>libc</code> , <code>collctype</code> is NULL and <code>colllocale</code> is used instead.
<code>colllocale text</code>	Collation provider locale name for this collation object. If the provider is <code>libc</code> , <code>colllocale</code> is NULL; <code>collcollate</code> and <code>collctype</code> are used instead.
<code>collicurules text</code>	ICU collation rules for this collation object
<code>collversion text</code>	Provider-specific version of the collation. This is recorded when the collation is created and then checked when it is used, to detect changes in the collation definition that could lead to data corruption.

Note that the unique key on this catalog is (`collname`, `collencoding`, `collnamespace`) not just (`collname`, `collnamespace`). Postgres Pro generally ignores all collations that do not have `collencoding` equal to either the current database's encoding or -1, and creation of new entries with the same name as an entry with `collencoding` = -1 is forbidden. Therefore it is sufficient to use a qualified SQL name (*schema.name*) to identify a collation, even though this is not unique according to the catalog definition. The reason for defining the catalog this way is that `initdb` fills it in at cluster initialization time with entries for all locales available on the system, so it must be able to hold entries for all encodings that might ever be used in the cluster.

In the `template0` database, it could be useful to create collations whose encoding does not match the database encoding, since they could match the encodings of databases later cloned from `template0`. This would currently have to be done manually.

54.13. `pg_constraint`

The catalog `pg_constraint` stores check, primary key, unique, foreign key, and exclusion constraints on tables, as well as not-null constraints on domains. (Column constraints are not treated specially. Every column constraint is equivalent to some table constraint.) Not-null constraints on relations are represented in the `pg_attribute` catalog, not here.

User-defined constraint triggers (created with `CREATE CONSTRAINT TRIGGER`) also give rise to an entry in this table.

Check constraints on domains are stored here, too.

Table 54.13. `pg_constraint` Columns

Column	Type	Description
<code>oid</code>	<code>oid</code>	Row identifier
<code>conname</code>	<code>name</code>	Constraint name (not necessarily unique!)
<code>connamespace</code>	<code>oid</code> (references <code>pg_namespace.oid</code>)	The OID of the namespace that contains this constraint
<code>contype</code>	<code>char</code>	<code>c</code> = check constraint, <code>f</code> = foreign key constraint, <code>n</code> = not-null constraint (domains only), <code>p</code> = primary key constraint, <code>u</code> = unique constraint, <code>t</code> = constraint trigger, <code>x</code> = exclusion constraint
<code>condeferrable</code>	<code>bool</code>	Is the constraint deferrable?
<code>condeferred</code>	<code>bool</code>	Is the constraint deferred by default?
<code>convalidated</code>	<code>bool</code>	Has the constraint been validated? Currently, can be false only for foreign keys and CHECK constraints
<code>conrelid</code>	<code>oid</code> (references <code>pg_class.oid</code>)	The table this constraint is on; zero if not a table constraint
<code>contypid</code>	<code>oid</code> (references <code>pg_type.oid</code>)	The domain this constraint is on; zero if not a domain constraint
<code>conindid</code>	<code>oid</code> (references <code>pg_class.oid</code>)	The index supporting this constraint, if it's a unique, primary key, foreign key, or exclusion constraint; else zero
<code>conparentid</code>	<code>oid</code> (references <code>pg_constraint.oid</code>)	The corresponding constraint of the parent partitioned table, if this is a constraint on a partition; else zero
<code>confrelid</code>	<code>oid</code> (references <code>pg_class.oid</code>)	If a foreign key, the referenced table; else zero
<code>confupdtype</code>	<code>char</code>	Foreign key update action code: <code>a</code> = no action, <code>r</code> = restrict, <code>c</code> = cascade, <code>n</code> = set null, <code>d</code> = set default
<code>confdeltype</code>	<code>char</code>	Foreign key deletion action code: <code>a</code> = no action, <code>r</code> = restrict, <code>c</code> = cascade, <code>n</code> = set null, <code>d</code> = set default
<code>conmatchtype</code>	<code>char</code>	Foreign key match type: <code>f</code> = full, <code>p</code> = partial, <code>s</code> = simple
<code>conislocal</code>	<code>bool</code>	This constraint is defined locally for the relation. Note that a constraint can be locally defined and inherited simultaneously.
<code>coninhcount</code>	<code>int2</code>	The number of direct inheritance ancestors this constraint has. A constraint with a nonzero number of ancestors cannot be dropped nor renamed.

Column Type	Description
<code>connoinherit bool</code>	This constraint is defined locally for the relation. It is a non-inheritable constraint.
<code>conkey int2[] (references pg_attribute .attnum)</code>	If a table constraint (including foreign keys, but not constraint triggers), list of the constrained columns
<code>confkey int2[] (references pg_attribute .attnum)</code>	If a foreign key, list of the referenced columns
<code>conpfeqop oid[] (references pg_operator .oid)</code>	If a foreign key, list of the equality operators for PK = FK comparisons
<code>conppeqop oid[] (references pg_operator .oid)</code>	If a foreign key, list of the equality operators for PK = PK comparisons
<code>conffeqop oid[] (references pg_operator .oid)</code>	If a foreign key, list of the equality operators for FK = FK comparisons
<code>confdelsetcols int2[] (references pg_attribute .attnum)</code>	If a foreign key with a SET NULL or SET DEFAULT delete action, the columns that will be updated. If null, all of the referencing columns will be updated.
<code>conexcllop oid[] (references pg_operator .oid)</code>	If an exclusion constraint, list of the per-column exclusion operators
<code>conbin pg_node_tree</code>	If a check constraint, an internal representation of the expression. (It's recommended to use <code>pg_get_constraintdef()</code> to extract the definition of a check constraint.)

In the case of an exclusion constraint, `conkey` is only useful for constraint elements that are simple column references. For other cases, a zero appears in `conkey` and the associated index must be consulted to discover the expression that is constrained. (`conkey` thus has the same contents as `pg_index.indkey` for the index.)

Note

`pg_class.relchecks` needs to agree with the number of check-constraint entries found in this table for each relation.

54.14. `pg_conversion`

The catalog `pg_conversion` describes encoding conversion functions. See [CREATE CONVERSION](#) for more information.

Table 54.14. `pg_conversion` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>conname name</code>	Conversion name (unique within a namespace)
<code>connamespace oid (references pg_namespace .oid)</code>	The OID of the namespace that contains this conversion
<code>conowner oid (references pg_authid .oid)</code>	

Column Type	Description
	Owner of the conversion
conforencoding int4	Source encoding ID (<code>pg_encoding_to_char()</code> can translate this number to the encoding name)
contoencoding int4	Destination encoding ID (<code>pg_encoding_to_char()</code> can translate this number to the encoding name)
conproc regproc (references <code>pg_proc</code> .oid)	Conversion function
condefault bool	True if this is the default conversion

54.15. pg_database

The catalog `pg_database` stores information about the available databases. Databases are created with the `CREATE DATABASE` command. Consult [Chapter 22](#) for details about the meaning of some of the parameters.

Unlike most system catalogs, `pg_database` is shared across all databases of a cluster: there is only one copy of `pg_database` per cluster, not one per database.

Table 54.15. pg_database Columns

Column Type	Description
oid oid	Row identifier
datname name	Database name
datdba oid (references <code>pg_authid</code> .oid)	Owner of the database, usually the user who created it
encoding int4	Character encoding for this database (<code>pg_encoding_to_char()</code> can translate this number to the encoding name)
datlocprovider char	Locale provider for this database: <code>b</code> = builtin, <code>c</code> = libc, <code>i</code> = icu
datistemplate bool	If true, then this database can be cloned by any user with <code>CREATEDB</code> privileges; if false, then only superusers or the owner of the database can clone it.
dataallowconn bool	If false then no one can connect to this database. This is used to protect the <code>template0</code> database from being altered.
dathasloginevt bool	Indicates that there are login event triggers defined for this database. This flag is used to avoid extra lookups on the <code>pg_event_trigger</code> table during each backend startup. This flag is used internally by Postgres Pro and should not be manually altered or read for monitoring purposes.
datconlimit int4	Sets maximum number of concurrent connections that can be made to this database. -1 means no limit, -2 indicates the database is invalid.

Column Type	Description
<code>datfrozenxid xid</code>	All transaction IDs before this one have been replaced with a permanent (“frozen”) transaction ID in this database. This is used to track whether the database needs to be vacuumed in order to allow <code>pg_xact</code> to be shrunk. It is the minimum of the per-table <code>pg_class .relfrozenxid</code> values.
<code>datminmxid xid</code>	All multixact IDs before this one have been replaced with a transaction ID in this database. This is used to track whether the database needs to be vacuumed in order to allow <code>pg_multixact</code> to be shrunk. It is the minimum of the per-table <code>pg_class .relminmxid</code> values.
<code>dattablespace oid (references <code>pg_tablespace .oid</code>)</code>	The default tablespace for the database. Within this database, all tables for which <code>pg_class .reltablespace</code> is zero will be stored in this tablespace; in particular, all the non-shared system catalogs will be there.
<code>datcollate text</code>	<code>LC_COLLATE</code> for this database
<code>datctype text</code>	<code>LC_CTYPE</code> for this database
<code>datlocale text</code>	Collation provider locale name for this database. If the provider is <code>libc</code> , <code>datlocale</code> is <code>NULL</code> ; <code>datcollate</code> and <code>datctype</code> are used instead.
<code>datcurules text</code>	ICU collation rules for this database
<code>datcollversion text</code>	Provider-specific version of the collation. This is recorded when the database is created and then checked when it is used, to detect changes in the collation definition that could lead to data corruption.
<code>datacl aclitem[]</code>	Access privileges; see Section 5.8 for details

54.16. `pg_db_role_setting`

The catalog `pg_db_role_setting` records the default values that have been set for run-time configuration variables, for each role and database combination.

Unlike most system catalogs, `pg_db_role_setting` is shared across all databases of a cluster: there is only one copy of `pg_db_role_setting` per cluster, not one per database.

Table 54.16. `pg_db_role_setting` Columns

Column Type	Description
<code>setdatabase oid (references <code>pg_database .oid</code>)</code>	The OID of the database the setting is applicable to, or zero if not database-specific
<code>setrole oid (references <code>pg_authid .oid</code>)</code>	The OID of the role the setting is applicable to, or zero if not role-specific
<code>setconfig text[]</code>	Defaults for run-time configuration variables

54.17. `pg_default_acl`

The catalog `pg_default_acl` stores initial privileges to be assigned to newly created objects.

Table 54.17. pg_default_acl Columns

Column Type	Description
oid oid	Row identifier
defaclrole oid (references pg_authid .oid)	The OID of the role associated with this entry
defaclnamespace oid (references pg_namespace .oid)	The OID of the namespace associated with this entry, or zero if none
defaclobjtype char	Type of object this entry is for: <i>r</i> = relation (table, view), <i>s</i> = sequence, <i>f</i> = function, <i>T</i> = type, <i>n</i> = schema
defaclacl aclitem[]	Access privileges that this type of object should have on creation

A `pg_default_acl` entry shows the initial privileges to be assigned to an object belonging to the indicated user. There are currently two types of entry: “global” entries with `defaclnamespace = zero`, and “per-schema” entries that reference a particular schema. If a global entry is present then it *overrides* the normal hard-wired default privileges for the object type. A per-schema entry, if present, represents privileges to be *added to* the global or hard-wired default privileges.

Note that when an ACL entry in another catalog is null, it is taken to represent the hard-wired default privileges for its object, *not* whatever might be in `pg_default_acl` at the moment. `pg_default_acl` is only consulted during object creation.

54.18. pg_depend

The catalog `pg_depend` records the dependency relationships between database objects. This information allows `DROP` commands to find which other objects must be dropped by `DROP CASCADE` or prevent dropping in the `DROP RESTRICT` case.

See also [pg_shdepend](#), which performs a similar function for dependencies involving objects that are shared across a database cluster.

Table 54.18. pg_depend Columns

Column Type	Description
classid oid (references pg_class .oid)	The OID of the system catalog the dependent object is in
objid oid (references any OID column)	The OID of the specific dependent object
objsubid int4	For a table column, this is the column number (the <code>objid</code> and <code>classid</code> refer to the table itself). For all other object types, this column is zero.
refclassid oid (references pg_class .oid)	The OID of the system catalog the referenced object is in
refobjid oid (references any OID column)	The OID of the specific referenced object
refobjsubid int4	For a table column, this is the column number (the <code>refobjid</code> and <code>refclassid</code> refer to the table itself). For all other object types, this column is zero.
deptype char	A code defining the specific semantics of this dependency relationship; see text

In all cases, a `pg_depend` entry indicates that the referenced object cannot be dropped without also dropping the dependent object. However, there are several subflavors identified by `deptype`:

`DEPENDENCY_NORMAL (n)`

A normal relationship between separately-created objects. The dependent object can be dropped without affecting the referenced object. The referenced object can only be dropped by specifying `CASCADE`, in which case the dependent object is dropped, too. Example: a table column has a normal dependency on its data type.

`DEPENDENCY_AUTO (a)`

The dependent object can be dropped separately from the referenced object, and should be automatically dropped (regardless of `RESTRICT` or `CASCADE` mode) if the referenced object is dropped. Example: a named constraint on a table is made auto-dependent on the table, so that it will go away if the table is dropped.

`DEPENDENCY_INTERNAL (i)`

The dependent object was created as part of creation of the referenced object, and is really just a part of its internal implementation. A direct `DROP` of the dependent object will be disallowed outright (we'll tell the user to issue a `DROP` against the referenced object, instead). A `DROP` of the referenced object will result in automatically dropping the dependent object whether `CASCADE` is specified or not. If the dependent object has to be dropped due to a dependency on some other object being removed, its drop is converted to a drop of the referenced object, so that `NORMAL` and `AUTO` dependencies of the dependent object behave much like they were dependencies of the referenced object. Example: a view's `ON SELECT` rule is made internally dependent on the view, preventing it from being dropped while the view remains. Dependencies of the rule (such as tables it refers to) act as if they were dependencies of the view.

`DEPENDENCY_PARTITION_PRI (P)`

`DEPENDENCY_PARTITION_SEC (S)`

The dependent object was created as part of creation of the referenced object, and is really just a part of its internal implementation; however, unlike `INTERNAL`, there is more than one such referenced object. The dependent object must not be dropped unless at least one of these referenced objects is dropped; if any one is, the dependent object should be dropped whether or not `CASCADE` is specified. Also unlike `INTERNAL`, a drop of some other object that the dependent object depends on does not result in automatic deletion of any partition-referenced object. Hence, if the drop does not cascade to at least one of these objects via some other path, it will be refused. (In most cases, the dependent object shares all its non-partition dependencies with at least one partition-referenced object, so that this restriction does not result in blocking any cascaded delete.) Primary and secondary partition dependencies behave identically except that the primary dependency is preferred for use in error messages; hence, a partition-dependent object should have one primary partition dependency and one or more secondary partition dependencies. Note that partition dependencies are made in addition to, not instead of, any dependencies the object would normally have. This simplifies `ATTACH/DETACH PARTITION` operations: the partition dependencies need only be added or removed. Example: a child partitioned index is made partition-dependent on both the partition table it is on and the parent partitioned index, so that it goes away if either of those is dropped, but not otherwise. The dependency on the parent index is primary, so that if the user tries to drop the child partitioned index, the error message will suggest dropping the parent index instead (not the table).

`DEPENDENCY_EXTENSION (e)`

The dependent object is a member of the *extension* that is the referenced object (see `pg_extension`). The dependent object can be dropped only via `DROP EXTENSION` on the referenced object. Functionally this dependency type acts the same as an `INTERNAL` dependency, but it's kept separate for clarity and to simplify `pg_dump`.

`DEPENDENCY_AUTO_EXTENSION (x)`

The dependent object is not a member of the extension that is the referenced object (and so it should not be ignored by `pg_dump`), but it cannot function without the extension and should be auto-dropped

if the extension is. The dependent object may be dropped on its own as well. Functionally this dependency type acts the same as an `AUTO` dependency, but it's kept separate for clarity and to simplify `pg_dump`.

Other dependency flavors might be needed in future.

Note that it's quite possible for two objects to be linked by more than one `pg_depend` entry. For example, a child partitioned index would have both a partition-type dependency on its associated partition table, and an auto dependency on each column of that table that it indexes. This sort of situation expresses the union of multiple dependency semantics. A dependent object can be dropped without `CASCADE` if any of its dependencies satisfies its condition for automatic dropping. Conversely, all the dependencies' restrictions about which objects must be dropped together must be satisfied.

Most objects created during `initdb` are considered “pinned”, which means that the system itself depends on them. Therefore, they are never allowed to be dropped. Also, knowing that pinned objects will not be dropped, the dependency mechanism doesn't bother to make `pg_depend` entries showing dependencies on them. Thus, for example, a table column of type `numeric` notionally has a `NORMAL` dependency on the `numeric` data type, but no such entry actually appears in `pg_depend`.

54.19. `pg_description`

The catalog `pg_description` stores optional descriptions (comments) for each database object. Descriptions can be manipulated with the `COMMENT` command and viewed with `psql`'s `\d` commands. Descriptions of many built-in system objects are provided in the initial contents of `pg_description`.

See also `pg_shdescription`, which performs a similar function for descriptions involving objects that are shared across a database cluster.

Table 54.19. `pg_description` Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	The OID of the object this description pertains to
<code>classoid oid</code> (references <code>pg_class</code> .oid)	The OID of the system catalog this object appears in
<code>objsubid int4</code>	For a comment on a table column, this is the column number (the <code>objoid</code> and <code>classoid</code> refer to the table itself). For all other object types, this column is zero.
<code>description text</code>	Arbitrary text that serves as the description of this object

54.20. `pg_enum`

The `pg_enum` catalog contains entries showing the values and labels for each enum type. The internal representation of a given enum value is actually the OID of its associated row in `pg_enum`.

Table 54.20. `pg_enum` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>enumtypid oid</code> (references <code>pg_type</code> .oid)	The OID of the <code>pg_type</code> entry owning this enum value
<code>enumsortorder float4</code>	

Column Type	Description
	The sort position of this enum value within its enum type
enumlabel name	The textual label for this enum value

The OIDs for `pg_enum` rows follow a special rule: even-numbered OIDs are guaranteed to be ordered in the same way as the sort ordering of their enum type. That is, if two even OIDs belong to the same enum type, the smaller OID must have the smaller `enumsortorder` value. Odd-numbered OID values need bear no relationship to the sort order. This rule allows the enum comparison routines to avoid catalog lookups in many common cases. The routines that create and alter enum types attempt to assign even OIDs to enum values whenever possible.

When an enum type is created, its members are assigned sort-order positions 1..*n*. But members added later might be given negative or fractional values of `enumsortorder`. The only requirement on these values is that they be correctly ordered and unique within each enum type.

54.21. pg_event_trigger

The catalog `pg_event_trigger` stores event triggers. See [Chapter 41](#) for more information.

Table 54.21. pg_event_trigger Columns

Column Type	Description
oid oid	Row identifier
evtname name	Trigger name (must be unique)
evtevent name	Identifies the event for which this trigger fires
evtowner oid (references pg_authid .oid)	Owner of the event trigger
evtfoid oid (references pg_proc .oid)	The function to be called
evtenabled char	Controls in which session_replication_role modes the event trigger fires. <code>O</code> = trigger fires in “origin” and “local” modes, <code>D</code> = trigger is disabled, <code>R</code> = trigger fires in “replica” mode, <code>A</code> = trigger fires always.
evttags text[]	Command tags for which this trigger will fire. If NULL, the firing of this trigger is not restricted on the basis of the command tag.

54.22. pg_extension

The catalog `pg_extension` stores information about the installed extensions. See [Section 39.17](#) for details about extensions.

Table 54.22. pg_extension Columns

Column Type	Description
oid oid	Row identifier
extname name	

Column Type	Description
	Name of the extension
extowner oid (references pg_authid .oid)	Owner of the extension
extnamespace oid (references pg_namespace .oid)	Schema containing the extension's exported objects
extrelocatable bool	True if extension can be relocated to another schema
extversion text	Version name for the extension
extconfig oid[] (references pg_class .oid)	Array of regclass OIDs for the extension's configuration table(s), or NULL if none
extcondition text[]	Array of WHERE-clause filter conditions for the extension's configuration table(s), or NULL if none

Note that unlike most catalogs with a “namespace” column, `extnamespace` is not meant to imply that the extension belongs to that schema. Extension names are never schema-qualified. Rather, `extnamespace` indicates the schema that contains most or all of the extension's objects. If `extrelocatable` is true, then this schema must in fact contain all schema-qualifiable objects belonging to the extension.

54.23. `pg_foreign_data_wrapper`

The catalog `pg_foreign_data_wrapper` stores foreign-data wrapper definitions. A foreign-data wrapper is the mechanism by which external data, residing on foreign servers, is accessed.

Table 54.23. `pg_foreign_data_wrapper` Columns

Column Type	Description
oid oid	Row identifier
fdwname name	Name of the foreign-data wrapper
fdwowner oid (references pg_authid .oid)	Owner of the foreign-data wrapper
fdwhandler oid (references pg_proc .oid)	References a handler function that is responsible for supplying execution routines for the foreign-data wrapper. Zero if no handler is provided
fdwvalidator oid (references pg_proc .oid)	References a validator function that is responsible for checking the validity of the options given to the foreign-data wrapper, as well as options for foreign servers and user mappings using the foreign-data wrapper. Zero if no validator is provided
fdwacl aclitem[]	Access privileges; see Section 5.8 for details
fdwoptions text[]	Foreign-data wrapper specific options, as “keyword=value” strings

54.24. `pg_foreign_server`

The catalog `pg_foreign_server` stores foreign server definitions. A foreign server describes a source of external data, such as a remote server. Foreign servers are accessed via foreign-data wrappers.

Table 54.24. pg_foreign_server Columns

Column Type	Description
oid oid	Row identifier
srvname name	Name of the foreign server
srvowner oid (references pg_authid .oid)	Owner of the foreign server
srvfdw oid (references pg_foreign_data_wrapper .oid)	OID of the foreign-data wrapper of this foreign server
srvtype text	Type of the server (optional)
srvversion text	Version of the server (optional)
srvacl aclitem[]	Access privileges; see Section 5.8 for details
srvoptions text[]	Foreign server specific options, as “keyword=value” strings

54.25. pg_foreign_table

The catalog `pg_foreign_table` contains auxiliary information about foreign tables. A foreign table is primarily represented by a [pg_class](#) entry, just like a regular table. Its `pg_foreign_table` entry contains the information that is pertinent only to foreign tables and not any other kind of relation.

Table 54.25. pg_foreign_table Columns

Column Type	Description
ftrelid oid (references pg_class .oid)	The OID of the pg_class entry for this foreign table
ftserver oid (references pg_foreign_server .oid)	OID of the foreign server for this foreign table
ftoptions text[]	Foreign table options, as “keyword=value” strings

54.26. pg_index

The catalog `pg_index` contains part of the information about indexes. The rest is mostly in [pg_class](#).

Table 54.26. pg_index Columns

Column Type	Description
indexrelid oid (references pg_class .oid)	The OID of the pg_class entry for this index
indrelid oid (references pg_class .oid)	The OID of the pg_class entry for the table this index is for
indnatts int2	The total number of columns in the index (duplicates <code>pg_class.relnatts</code>); this number includes both key and included attributes
indnkeyatts int2	

Column Type	Description
	The number of <i>key columns</i> in the index, not counting any <i>included columns</i> , which are merely stored and do not participate in the index semantics
<code>indisunique</code> bool	If true, this is a unique index
<code>indnullsnotdistinct</code> bool	This value is only used for unique indexes. If false, this unique index will consider null values distinct (so the index can contain multiple null values in a column, the default Postgres Pro behavior). If it is true, it will consider null values to be equal (so the index can only contain one null value in a column).
<code>indisprimary</code> bool	If true, this index represents the primary key of the table (<code>indisunique</code> should always be true when this is true)
<code>indisexclusion</code> bool	If true, this index supports an exclusion constraint
<code>indimmediate</code> bool	If true, the uniqueness check is enforced immediately on insertion (irrelevant if <code>indisunique</code> is not true)
<code>indisclustered</code> bool	If true, the table was last clustered on this index
<code>indisvalid</code> bool	If true, the index is currently valid for queries. False means the index is possibly incomplete: it must still be modified by <code>INSERT/UPDATE</code> operations, but it cannot safely be used for queries. If it is unique, the uniqueness property is not guaranteed true either.
<code>indcheckxmin</code> bool	If true, queries must not use the index until the <code>xmin</code> of this <code>pg_index</code> row is below their <code>TransactionXmin</code> event horizon, because the table may contain broken HOT chains with incompatible rows that they can see
<code>indisready</code> bool	If true, the index is currently ready for inserts. False means the index must be ignored by <code>INSERT/UPDATE</code> operations.
<code>indislive</code> bool	If false, the index is in process of being dropped, and should be ignored for all purposes (including HOT-safety decisions)
<code>indisreplident</code> bool	If true this index has been chosen as “replica identity” using <code>ALTER TABLE ... REPLICA IDENTITY USING INDEX ...</code>
<code>indkey</code> int2vector (references <code>pg_attribute.attnum</code>)	This is an array of <code>indnatts</code> values that indicate which table columns this index indexes. For example, a value of 1 3 would mean that the first and the third table columns make up the index entries. Key columns come before non-key (included) columns. A zero in this array indicates that the corresponding index attribute is an expression over the table columns, rather than a simple column reference.
<code>indcollation</code> oidvector (references <code>pg_collation.oid</code>)	For each column in the index key (<code>indnkeyatts</code> values), this contains the OID of the collation to use for the index, or zero if the column is not of a collatable data type.
<code>indclass</code> oidvector (references <code>pg_opclass.oid</code>)	For each column in the index key (<code>indnkeyatts</code> values), this contains the OID of the operator class to use. See <code>pg_opclass</code> for details.
<code>indoption</code> int2vector	

Column Type	Description
	This is an array of <code>indnkeyatts</code> values that store per-column flag bits. The meaning of the bits is defined by the index's access method.
<code>indexprs pg_node_tree</code>	Expression trees (in <code>nodeToString()</code> representation) for index attributes that are not simple column references. This is a list with one element for each zero entry in <code>indkey</code> . Null if all index attributes are simple references.
<code>indpred pg_node_tree</code>	Expression tree (in <code>nodeToString()</code> representation) for partial index predicate. Null if not a partial index.

54.27. `pg_inherits`

The catalog `pg_inherits` records information about table and index inheritance hierarchies. There is one entry for each direct parent-child table or index relationship in the database. (Indirect inheritance can be determined by following chains of entries.)

Table 54.27. `pg_inherits` Columns

Column Type	Description
<code>inhrelid oid</code> (references <code>pg_class .oid</code>)	The OID of the child table or index
<code>inhparent oid</code> (references <code>pg_class .oid</code>)	The OID of the parent table or index
<code>inhseqno int4</code>	If there is more than one direct parent for a child table (multiple inheritance), this number tells the order in which the inherited columns are to be arranged. The count starts at 1. Indexes cannot have multiple inheritance, since they can only inherit when using declarative partitioning.
<code>inhdetachpending bool</code>	<code>true</code> for a partition that is in the process of being detached; <code>false</code> otherwise.

54.28. `pg_init_privs`

The catalog `pg_init_privs` records information about the initial privileges of objects in the system. There is one entry for each object in the database which has a non-default (non-NULL) initial set of privileges.

Objects can have initial privileges either by having those privileges set when the system is initialized (by `initdb`) or when the object is created during a `CREATE EXTENSION` and the extension script sets initial privileges using the `GRANT` system. Note that the system will automatically handle recording of the privileges during the extension script and that extension authors need only use the `GRANT` and `REVOKE` statements in their script to have the privileges recorded. The `privtype` column indicates if the initial privilege was set by `initdb` or during a `CREATE EXTENSION` command.

Objects which have initial privileges set by `initdb` will have entries where `privtype` is 'i', while objects which have initial privileges set by `CREATE EXTENSION` will have entries where `privtype` is 'e'.

Table 54.28. `pg_init_privs` Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	

Column Type	Description
	The OID of the specific object
<code>classoid oid</code> (references <code>pg_class</code> .oid)	The OID of the system catalog the object is in
<code>objsubid int4</code>	For a table column, this is the column number (the <code>objoid</code> and <code>classoid</code> refer to the table itself). For all other object types, this column is zero.
<code>privtype char</code>	A code defining the type of initial privilege of this object; see text
<code>initprivs aclitem[]</code>	The initial access privileges; see Section 5.8 for details

54.29. pg_language

The catalog `pg_language` registers languages in which you can write functions or stored procedures. See [CREATE LANGUAGE](#) and [Chapter 43](#) for more information about language handlers.

Table 54.29. pg_language Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>lanname name</code>	Name of the language
<code>lanowner oid</code> (references <code>pg_authid</code> .oid)	Owner of the language
<code>lanispl bool</code>	This is false for internal languages (such as SQL) and true for user-defined languages. Currently, <code>pg_dump</code> still uses this to determine which languages need to be dumped, but this might be replaced by a different mechanism in the future.
<code>lanpltrusted bool</code>	True if this is a trusted language, which means that it is believed not to grant access to anything outside the normal SQL execution environment. Only superusers can create functions in untrusted languages.
<code>lanplcallfoid oid</code> (references <code>pg_proc</code> .oid)	For noninternal languages this references the language handler, which is a special function that is responsible for executing all functions that are written in the particular language. Zero for internal languages.
<code>laninline oid</code> (references <code>pg_proc</code> .oid)	This references a function that is responsible for executing “inline” anonymous code blocks (DO blocks). Zero if inline blocks are not supported.
<code>lanvalidator oid</code> (references <code>pg_proc</code> .oid)	This references a language validator function that is responsible for checking the syntax and validity of new functions when they are created. Zero if no validator is provided.
<code>lanacl aclitem[]</code>	Access privileges; see Section 5.8 for details

54.30. pg_largeobject

The catalog `pg_largeobject` holds the data making up “large objects”. A large object is identified by an OID assigned when it is created. Each large object is broken into segments or “pages” small enough

to be conveniently stored as rows in `pg_largeobject`. The amount of data per page is defined to be `LOBLKSIZE` (which is currently `BLCKSZ/4`, or typically 2 kB).

Prior to PostgreSQL 9.0, there was no permission structure associated with large objects. As a result, `pg_largeobject` was publicly readable and could be used to obtain the OIDs (and contents) of all large objects in the system. This is no longer the case; use `pg_largeobject_metadata` to obtain a list of large object OIDs.

Table 54.30. `pg_largeobject` Columns

Column Type	Description
<code>loid oid</code> (references <code>pg_largeobject_metadata .oid</code>)	Identifier of the large object that includes this page
<code>pageno int4</code>	Page number of this page within its large object (counting from zero)
<code>data bytea</code>	Actual data stored in the large object. This will never be more than <code>LOBLKSIZE</code> bytes and might be less.

Each row of `pg_largeobject` holds data for one page of a large object, beginning at byte offset (`pageno * LOBLKSIZE`) within the object. The implementation allows sparse storage: pages might be missing, and might be shorter than `LOBLKSIZE` bytes even if they are not the last page of the object. Missing regions within a large object read as zeroes.

54.31. `pg_largeobject_metadata`

The catalog `pg_largeobject_metadata` holds metadata associated with large objects. The actual large object data is stored in `pg_largeobject`.

Table 54.31. `pg_largeobject_metadata` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>lomowner oid</code> (references <code>pg_authid .oid</code>)	Owner of the large object
<code>lomacl aclitem[]</code>	Access privileges; see Section 5.8 for details

54.32. `pg_namespace`

The catalog `pg_namespace` stores namespaces. A namespace is the structure underlying SQL schemas: each namespace can have a separate collection of relations, types, etc. without name conflicts.

Table 54.32. `pg_namespace` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>nspname name</code>	Name of the namespace
<code>nspowner oid</code> (references <code>pg_authid .oid</code>)	Owner of the namespace
<code>nspsecofficer oid</code>	

Column Type	Description
	Security officer of the namespace
nspacl aclitem[]	Access privileges; see Section 5.8 for details

54.33. pg_opclass

The catalog `pg_opclass` defines index access method operator classes. Each operator class defines semantics for index columns of a particular data type and a particular index access method. An operator class essentially specifies that a particular operator family is applicable to a particular indexable column data type. The set of operators from the family that are actually usable with the indexed column are whichever ones accept the column's data type as their left-hand input.

Operator classes are described at length in [Section 39.16](#).

Table 54.33. pg_opclass Columns

Column Type	Description
oid oid	Row identifier
opcmethod oid (references pg_am .oid)	Index access method operator class is for
opcname name	Name of this operator class
opcnamespace oid (references pg_namespace .oid)	Namespace of this operator class
opcowner oid (references pg_authid .oid)	Owner of the operator class
opcfamily oid (references pg_opfamily .oid)	Operator family containing the operator class
opcintype oid (references pg_type .oid)	Data type that the operator class indexes
opcdefault bool	True if this operator class is the default for opcintype
opckeytype oid (references pg_type .oid)	Type of data stored in index, or zero if same as opcintype

An operator class's `opcmethod` must match the `opfmeth` of its containing operator family. Also, there must be no more than one `pg_opclass` row having `opcdefault` true for any given combination of `opcmethod` and `opcintype`.

54.34. pg_operator

The catalog `pg_operator` stores information about operators. See [CREATE OPERATOR](#) and [Section 39.14](#) for more information.

Table 54.34. pg_operator Columns

Column Type	Description
oid oid	Row identifier
oprname name	

Column Type	Description
	Name of the operator
oprnamespace oid (references pg_namespace .oid)	The OID of the namespace that contains this operator
oprowner oid (references pg_authid .oid)	Owner of the operator
oprkind char	b = infix operator (“both”), or l = prefix operator (“left”)
oprcanmerge bool	This operator supports merge joins
oprcanhash bool	This operator supports hash joins
oprleft oid (references pg_type .oid)	Type of the left operand (zero for a prefix operator)
oprright oid (references pg_type .oid)	Type of the right operand
oprresult oid (references pg_type .oid)	Type of the result (zero for a not-yet-defined “shell” operator)
oprcom oid (references pg_operator .oid)	Commutator of this operator (zero if none)
oprnegate oid (references pg_operator .oid)	Negator of this operator (zero if none)
oprcode regproc (references pg_proc .oid)	Function that implements this operator (zero for a not-yet-defined “shell” operator)
oprrrest regproc (references pg_proc .oid)	Restriction selectivity estimation function for this operator (zero if none)
oprjoin regproc (references pg_proc .oid)	Join selectivity estimation function for this operator (zero if none)

54.35. pg_opfamily

The catalog `pg_opfamily` defines operator families. Each operator family is a collection of operators and associated support routines that implement the semantics specified for a particular index access method. Furthermore, the operators in a family are all “compatible”, in a way that is specified by the access method. The operator family concept allows cross-data-type operators to be used with indexes and to be reasoned about using knowledge of access method semantics.

Operator families are described at length in [Section 39.16](#).

Table 54.35. `pg_opfamily` Columns

Column Type	Description
oid oid	Row identifier
opfmethod oid (references pg_am .oid)	Index access method operator family is for
opfname name	Name of this operator family
opfnamespace oid (references pg_namespace .oid)	

Column Type	Description
	Namespace of this operator family
opfowner oid (references pg_authid .oid)	Owner of the operator family

The majority of the information defining an operator family is not in its `pg_opfamily` row, but in the associated rows in `pg_amop`, `pg_amproc`, and `pg_opclass`.

54.36. `pg_parameter_acl`

The catalog `pg_parameter_acl` records configuration parameters for which privileges have been granted to one or more roles. No entry is made for parameters that have default privileges.

Unlike most system catalogs, `pg_parameter_acl` is shared across all databases of a cluster: there is only one copy of `pg_parameter_acl` per cluster, not one per database.

Table 54.36. `pg_parameter_acl` Columns

Column Type	Description
oid oid	Row identifier
parname text	The name of a configuration parameter for which privileges are granted
paracl aclitem[]	Access privileges; see Section 5.8 for details

54.37. `pg_partitioned_table`

The catalog `pg_partitioned_table` stores information about how tables are partitioned.

Table 54.37. `pg_partitioned_table` Columns

Column Type	Description
partrelid oid (references pg_class .oid)	The OID of the pg_class entry for this partitioned table
partstrat char	Partitioning strategy; h = hash partitioned table, l = list partitioned table, r = range partitioned table
partnatts int2	The number of columns in the partition key
partdefid oid (references pg_class .oid)	The OID of the pg_class entry for the default partition of this partitioned table, or zero if this partitioned table does not have a default partition
partattrs int2vector (references pg_attribute .attnum)	This is an array of <code>partnatts</code> values that indicate which table columns are part of the partition key. For example, a value of 1 3 would mean that the first and the third table columns make up the partition key. A zero in this array indicates that the corresponding partition key column is an expression, rather than a simple column reference.
partclass oidvector (references pg_opclass .oid)	For each column in the partition key, this contains the OID of the operator class to use. See pg_opclass for details.
partcollation oidvector (references pg_collation .oid)	

Column Type	Description
	For each column in the partition key, this contains the OID of the collation to use for partitioning, or zero if the column is not of a collatable data type.
partexprs pg_node_tree	Expression trees (in <code>nodeToString()</code> representation) for partition key columns that are not simple column references. This is a list with one element for each zero entry in <code>partattrs</code> . Null if all partition key columns are simple references.

54.38. pg_policy

The catalog `pg_policy` stores row-level security policies for tables. A policy includes the kind of command that it applies to (possibly all commands), the roles that it applies to, the expression to be added as a security-barrier qualification to queries that include the table, and the expression to be added as a `WITH CHECK` option for queries that attempt to add new records to the table.

Table 54.38. pg_policy Columns

Column Type	Description
oid oid	Row identifier
polname name	The name of the policy
polrelid oid (references <code>pg_class</code> .oid)	The table to which the policy applies
polcmd char	The command type to which the policy is applied: <code>r</code> for <code>SELECT</code> , <code>a</code> for <code>INSERT</code> , <code>w</code> for <code>UPDATE</code> , <code>d</code> for <code>DELETE</code> , or <code>*</code> for all
polpermissive bool	Is the policy permissive or restrictive?
polroles oid[] (references <code>pg_authid</code> .oid)	The roles to which the policy is applied; zero means <code>PUBLIC</code> (and normally appears alone in the array)
polqual pg_node_tree	The expression tree to be added to the security barrier qualifications for queries that use the table
polwithcheck pg_node_tree	The expression tree to be added to the <code>WITH CHECK</code> qualifications for queries that attempt to add rows to the table

Note

Policies stored in `pg_policy` are applied only when `pg_class.relrowsecurity` is set for their table.

54.39. pg_proc

The catalog `pg_proc` stores information about functions, procedures, aggregate functions, and window functions (collectively also known as routines). See [CREATE FUNCTION](#), [CREATE PROCEDURE](#), and [Section 39.3](#) for more information.

If `prokind` indicates that the entry is for an aggregate function, there should be a matching row in `pg_aggregate`.

Table 54.39. pg_proc Columns

Column	Type	Description
oid	oid	Row identifier
proname	name	Name of the function
pronamespace	oid (references pg_namespace .oid)	The OID of the namespace that contains this function
proowner	oid (references pg_authid .oid)	Owner of the function
prolang	oid (references pg_language .oid)	Implementation language or call interface of this function
procost	float4	Estimated execution cost (in units of cpu_operator_cost); if <code>proretset</code> , this is cost per row returned
prorows	float4	Estimated number of result rows (zero if not <code>proretset</code>)
provariadic	oid (references pg_type .oid)	Data type of the variadic array parameter's elements, or zero if the function does not have a variadic parameter
prosupport	regproc (references pg_proc .oid)	Planner support function for this function (see Section 39.11), or zero if none
prokind	char	<code>f</code> for a normal function, <code>p</code> for a procedure, <code>a</code> for an aggregate function, or <code>w</code> for a window function
prosecdef	bool	Function is a security definer (i.e., a “setuid” function)
proleakproof	bool	The function has no side effects. No information about the arguments is conveyed except via the return value. Any function that might throw an error depending on the values of its arguments is not leak-proof.
proisstrict	bool	Function returns null if any call argument is null. In that case the function won't actually be called at all. Functions that are not “strict” must be prepared to handle null inputs.
proretset	bool	Function returns a set (i.e., multiple values of the specified data type)
provolatile	char	<code>provolatile</code> tells whether the function's result depends only on its input arguments, or is affected by outside factors. It is <code>i</code> for “immutable” functions, which always deliver the same result for the same inputs. It is <code>s</code> for “stable” functions, whose results (for fixed inputs) do not change within a scan. It is <code>v</code> for “volatile” functions, whose results might change at any time. (Use <code>v</code> also for functions with side-effects, so that calls to them cannot get optimized away.)
proparallel	char	<code>proparallel</code> tells whether the function can be safely run in parallel mode. It is <code>s</code> for functions which are safe to run in parallel mode without restriction. It is <code>r</code> for functions which can be run in parallel mode, but their execution is restricted to the parallel group leader; parallel worker processes cannot invoke these functions. It is <code>u</code> for functions which are unsafe in parallel mode; the presence of such a function forces a serial execution plan.
pronargs	int2	

Column Type	Description
	Number of input arguments
<code>pronargdefaults int2</code>	Number of arguments that have defaults
<code>prorettytype oid</code> (references pg_type .oid)	Data type of the return value
<code>proargtypes oidvector</code> (references pg_type .oid)	An array of the data types of the function arguments. This includes only input arguments (including <code>INOUT</code> and <code>VARIADIC</code> arguments), and thus represents the call signature of the function.
<code>proallargtypes oid[]</code> (references pg_type .oid)	An array of the data types of the function arguments. This includes all arguments (including <code>OUT</code> and <code>INOUT</code> arguments); however, if all the arguments are <code>IN</code> arguments, this field will be null. Note that subscripting is 1-based, whereas for historical reasons <code>proargtypes</code> is subscripted from 0.
<code>proargmodes char[]</code>	An array of the modes of the function arguments, encoded as <code>i</code> for <code>IN</code> arguments, <code>o</code> for <code>OUT</code> arguments, <code>b</code> for <code>INOUT</code> arguments, <code>v</code> for <code>VARIADIC</code> arguments, <code>t</code> for <code>TABLE</code> arguments. If all the arguments are <code>IN</code> arguments, this field will be null. Note that subscripts correspond to positions of <code>proallargtypes</code> not <code>proargtypes</code> .
<code>proargnames text[]</code>	An array of the names of the function arguments. Arguments without a name are set to empty strings in the array. If none of the arguments have a name, this field will be null. Note that subscripts correspond to positions of <code>proallargtypes</code> not <code>proargtypes</code> .
<code>proargdefaults pg_node_tree</code>	Expression trees (in <code>nodeToString()</code> representation) for default values. This is a list with <code>pronargdefaults</code> elements, corresponding to the last <i>N</i> input arguments (i.e., the last <i>N</i> <code>proargtypes</code> positions). If none of the arguments have defaults, this field will be null.
<code>protrftypes oid[]</code> (references pg_type .oid)	An array of the argument/result data type(s) for which to apply transforms (from the function's <code>TRANSFORM</code> clause). Null if none.
<code>prosrc text</code>	This tells the function handler how to invoke the function. It might be the actual source code of the function for interpreted languages, a link symbol, a file name, or just about anything else, depending on the implementation language/call convention.
<code>probin text</code>	Additional information about how to invoke the function. Again, the interpretation is language-specific.
<code>prosqlbody pg_node_tree</code>	Pre-parsed SQL function body. This is used for SQL-language functions when the body is given in SQL-standard notation rather than as a string literal. It's null in other cases.
<code>proconfig text[]</code>	Function's local settings for run-time configuration variables
<code>proacl aclitem[]</code>	Access privileges; see Section 5.8 for details

For compiled functions, both built-in and dynamically loaded, `prosrc` contains the function's C-language name (link symbol). For SQL-language functions, `prosrc` contains the function's source text if that is specified as a string literal; but if the function body is specified in SQL-standard style, `prosrc` is unused (typically it's an empty string) and `prosqlbody` contains the pre-parsed definition. For all other current-

ly-known language types, `prosrc` contains the function's source text. `probin` is null except for dynamically-loaded C functions, for which it gives the name of the shared library file containing the function.

54.40. `pg_profile`

The catalog `pg_profile` lists all profiles available for your cluster. For details on profiles, see [CREATE PROFILE](#).

Profiles are cluster-wide, so `pg_profile` is shared across all databases of a cluster. This catalog is not publicly readable.

Table 54.40. `pg_profile` Columns

Column Type	Description
oid oid	Row identifier
pflname name	Profile name
pflfailedloginattempts int4 FAILED_LOGIN_ATTEMPTS	parameter value
pflfailedauthkeepstime int4 FAILED_AUTH_KEEP_TIME	parameter value (in seconds)
pfluserinactivetime int4 USER_INACTIVE_TIME	parameter value (in seconds)
pflpasswordreusetime int4 PASSWORD_REUSE_TIME	parameter value (in seconds)
pflpasswordreusemax int4 PASSWORD_REUSE_MAX	parameter value
pflpasswordlifetime int4 PASSWORD_LIFE_TIME	parameter value (in seconds)
pflpasswordgracetime int4 PASSWORD_GRACE_TIME	parameter value (in seconds)
pflpasswordminuniqchars int4 PASSWORD_MIN_UNIQUE_CHARS	parameter value
pflpasswordminlen int4 PASSWORD_MIN_LEN	parameter value
pflpasswordrequirecomplex int4 PASSWORD_REQUIRE_COMPLEX	parameter value

The `UNLIMITED` and `DEFAULT` parameter values are stored in `pg_profile` as `-1` and `-2`, respectively.

54.41. `pg_publication`

The catalog `pg_publication` contains all publications created in the database. For more on publications see [Section 29.1](#).

Table 54.41. `pg_publication` Columns

Column Type	Description
oid oid	Row identifier
pubname name	

Column Type	Description
	Name of the publication
pubowner oid (references pg_authid .oid)	Owner of the publication
puballtables bool	If true, this publication automatically includes all tables in the database, including any that will be created in the future.
pubinsert bool	If true, INSERT operations are replicated for tables in the publication.
pubupdate bool	If true, UPDATE operations are replicated for tables in the publication.
pubdelete bool	If true, DELETE operations are replicated for tables in the publication.
pubtruncate bool	If true, TRUNCATE operations are replicated for tables in the publication.
pubviaroot bool	If true, operations on a leaf partition are replicated using the identity and schema of its top-most partitioned ancestor mentioned in the publication instead of its own.

54.42. pg_publication_namespace

The catalog `pg_publication_namespace` contains the mapping between schemas and publications in the database. This is a many-to-many mapping.

Table 54.42. pg_publication_namespace Columns

Column Type	Description
oid oid	Row identifier
pnpubid oid (references pg_publication .oid)	Reference to publication
pnnspid oid (references pg_namespace .oid)	Reference to schema

54.43. pg_publication_rel

The catalog `pg_publication_rel` contains the mapping between relations and publications in the database. This is a many-to-many mapping. See also [Section 55.18](#) for a more user-friendly view of this information.

Table 54.43. pg_publication_rel Columns

Column Type	Description
oid oid	Row identifier
prpubid oid (references pg_publication .oid)	Reference to publication
prrelid oid (references pg_class .oid)	Reference to relation
prqual pg_node_tree	

Column Type	Description
	Expression tree (in <code>nodeToString()</code> representation) for the relation's publication qualifying condition. Null if there is no publication qualifying condition.
<code>prattrs int2vector</code> (references <code>pg_attribute.attnum</code>)	This is an array of values that indicates which table columns are part of the publication. For example, a value of <code>1 3</code> would mean that the first and the third table columns are published. A null value indicates that all columns are published.

54.44. `pg_range`

The catalog `pg_range` stores information about range types. This is in addition to the types' entries in `pg_type`.

Table 54.44. `pg_range` Columns

Column Type	Description
<code>rngtypeid oid</code> (references <code>pg_type.oid</code>)	OID of the range type
<code>rngsubtype oid</code> (references <code>pg_type.oid</code>)	OID of the element type (subtype) of this range type
<code>rngmultitypid oid</code> (references <code>pg_type.oid</code>)	OID of the multirange type for this range type
<code>rngcollation oid</code> (references <code>pg_collation.oid</code>)	OID of the collation used for range comparisons, or zero if none
<code>rngsubopc oid</code> (references <code>pg_opclass.oid</code>)	OID of the subtype's operator class used for range comparisons
<code>rngcanonical regproc</code> (references <code>pg_proc.oid</code>)	OID of the function to convert a range value into canonical form, or zero if none
<code>rngsubdiff regproc</code> (references <code>pg_proc.oid</code>)	OID of the function to return the difference between two element values as double precision, or zero if none

`rngsubopc` (plus `rngcollation`, if the element type is collatable) determines the sort ordering used by the range type. `rngcanonical` is used when the element type is discrete. `rngsubdiff` is optional but should be supplied to improve performance of GiST indexes on the range type.

54.45. `pg_replication_origin`

The `pg_replication_origin` catalog contains all replication origins created. For more on replication origins see [Chapter 51](#).

Unlike most system catalogs, `pg_replication_origin` is shared across all databases of a cluster: there is only one copy of `pg_replication_origin` per cluster, not one per database.

Table 54.45. `pg_replication_origin` Columns

Column Type	Description
<code>roident oid</code>	A unique, cluster-wide identifier for the replication origin. Should never leave the system.
<code>roname text</code>	The external, user defined, name of a replication origin.

54.46. pg_rewrite

The catalog `pg_rewrite` stores rewrite rules for tables and views.

Table 54.46. pg_rewrite Columns

Column	Type	Description
<code>oid</code>	<code>oid</code>	Row identifier
<code>rulename</code>	<code>name</code>	Rule name
<code>ev_class</code>	<code>oid</code> (references <code>pg_class.oid</code>)	The table this rule is for
<code>ev_type</code>	<code>char</code>	Event type that the rule is for: 1 = <code>SELECT</code> , 2 = <code>UPDATE</code> , 3 = <code>INSERT</code> , 4 = <code>DELETE</code>
<code>ev_enabled</code>	<code>char</code>	Controls in which <code>session_replication_role</code> modes the rule fires. <code>O</code> = rule fires in “origin” and “local” modes, <code>D</code> = rule is disabled, <code>R</code> = rule fires in “replica” mode, <code>A</code> = rule fires always.
<code>is_instead</code>	<code>bool</code>	True if the rule is an <code>INSTEAD</code> rule
<code>ev_qual</code>	<code>pg_node_tree</code>	Expression tree (in the form of a <code>nodeToString()</code> representation) for the rule's qualifying condition
<code>ev_action</code>	<code>pg_node_tree</code>	Query tree (in the form of a <code>nodeToString()</code> representation) for the rule's action

Note

`pg_class.relhasrules` must be true if a table has any rules in this catalog.

54.47. pg_role_password

The catalog `pg_role_password` records role's passwords history used to enforce authentication restrictions.

Unlike most system catalogs, `pg_role_password` is shared across all databases of a cluster: there is only one copy of `pg_role_password` per cluster, not one per database.

Table 54.47. pg_role_password Columns

Column	Type	Description
<code>passsetat</code>	<code>timestamp_tz</code>	Password set time
<code>passrole</code>	<code>regprofile</code> (references <code>pg_authid.oid</code>)	The OID of the role
<code>passpassword</code>	<code>text</code>	Password of the role (possibly encrypted)

54.48. pg_seclabel

The catalog `pg_seclabel` stores security labels on database objects. Security labels can be manipulated with the `SECURITY LABEL` command. For an easier way to view security labels, see [Section 55.23](#).

See also `pg_shseclabel`, which performs a similar function for security labels of database objects that are shared across a database cluster.

Table 54.48. `pg_seclabel` Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	The OID of the object this security label pertains to
<code>classoid oid</code> (references <code>pg_class</code> .oid)	The OID of the system catalog this object appears in
<code>objsubid int4</code>	For a security label on a table column, this is the column number (the <code>objoid</code> and <code>classoid</code> refer to the table itself). For all other object types, this column is zero.
<code>provider text</code>	The label provider associated with this label.
<code>label text</code>	The security label applied to this object.

54.49. `pg_sequence`

The catalog `pg_sequence` contains information about sequences. Some of the information about sequences, such as the name and the schema, is in `pg_class`

Table 54.49. `pg_sequence` Columns

Column Type	Description
<code>seqrelid oid</code> (references <code>pg_class</code> .oid)	The OID of the <code>pg_class</code> entry for this sequence
<code>seqtypeid oid</code> (references <code>pg_type</code> .oid)	Data type of the sequence
<code>seqstart int8</code>	Start value of the sequence
<code>seqincrement int8</code>	Increment value of the sequence
<code>seqmax int8</code>	Maximum value of the sequence
<code>seqmin int8</code>	Minimum value of the sequence
<code>seqcache int8</code>	Cache size of the sequence
<code>seqcycle bool</code>	Whether the sequence cycles

54.50. `pg_shdepend`

The catalog `pg_shdepend` records the dependency relationships between database objects and shared objects, such as roles. This information allows Postgres Pro to ensure that those objects are unreferenced before attempting to delete them.

See also [pg_depend](#), which performs a similar function for dependencies involving objects within a single database.

Unlike most system catalogs, `pg_shdepend` is shared across all databases of a cluster: there is only one copy of `pg_shdepend` per cluster, not one per database.

Table 54.50. `pg_shdepend` Columns

Column Type	Description
<code>dbid oid</code> (references pg_database .oid)	The OID of the database the dependent object is in, or zero for a shared object
<code>classid oid</code> (references pg_class .oid)	The OID of the system catalog the dependent object is in
<code>objid oid</code> (references any OID column)	The OID of the specific dependent object
<code>objsubid int4</code>	For a table column, this is the column number (the <code>objid</code> and <code>classid</code> refer to the table itself). For all other object types, this column is zero.
<code>refclassid oid</code> (references pg_class .oid)	The OID of the system catalog the referenced object is in (must be a shared catalog)
<code>refobjid oid</code> (references any OID column)	The OID of the specific referenced object
<code>deptype char</code>	A code defining the specific semantics of this dependency relationship; see text

In all cases, a `pg_shdepend` entry indicates that the referenced object cannot be dropped without also dropping the dependent object. However, there are several subflavors identified by `deptype`:

`SHARED_DEPENDENCY_OWNER` (o)

The referenced object (which must be a role) is the owner of the dependent object.

`SHARED_DEPENDENCY_SECURITY_OFFICER` (s)

The referenced object (which must be a role) is the security officer of the dependent object.

`SHARED_DEPENDENCY_ACL` (a)

The referenced object (which must be a role) is mentioned in the ACL (access control list, i.e., privileges list) of the dependent object. (A `SHARED_DEPENDENCY_ACL` entry is not made for the owner of the object, since the owner will have a `SHARED_DEPENDENCY_OWNER` entry anyway.)

`SHARED_DEPENDENCY_INITACL` (i)

The referenced object (which must be a role) is mentioned in a [pg_init_privs](#) entry for the dependent object.

`SHARED_DEPENDENCY_POLICY` (r)

The referenced object (which must be a role) is mentioned as the target of a dependent policy object.

`SHARED_DEPENDENCY_TABLESPACE` (t)

The referenced object (which must be a tablespace) is mentioned as the tablespace for a relation that doesn't have storage.

Other dependency flavors might be needed in future. Note in particular that the current definition only supports roles and tablespaces as referenced objects.

As in the `pg_depend` catalog, most objects created during `initdb` are considered “pinned”. No entries are made in `pg_shdepend` that would have a pinned object as either referenced or dependent object.

54.51. `pg_shdescription`

The catalog `pg_shdescription` stores optional descriptions (comments) for shared database objects. Descriptions can be manipulated with the `COMMENT` command and viewed with `psql`'s `\d` commands.

See also `pg_description`, which performs a similar function for descriptions involving objects within a single database.

Unlike most system catalogs, `pg_shdescription` is shared across all databases of a cluster: there is only one copy of `pg_shdescription` per cluster, not one per database.

Table 54.51. `pg_shdescription` Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	The OID of the object this description pertains to
<code>classoid oid</code> (references <code>pg_class .oid</code>)	The OID of the system catalog this object appears in
<code>description text</code>	Arbitrary text that serves as the description of this object

54.52. `pg_shseclabel`

The catalog `pg_shseclabel` stores security labels on shared database objects. Security labels can be manipulated with the `SECURITY LABEL` command. For an easier way to view security labels, see [Section 55.23](#).

See also `pg_seclabel`, which performs a similar function for security labels involving objects within a single database.

Unlike most system catalogs, `pg_shseclabel` is shared across all databases of a cluster: there is only one copy of `pg_shseclabel` per cluster, not one per database.

Table 54.52. `pg_shseclabel` Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	The OID of the object this security label pertains to
<code>classoid oid</code> (references <code>pg_class .oid</code>)	The OID of the system catalog this object appears in
<code>provider text</code>	The label provider associated with this label.
<code>label text</code>	The security label applied to this object.

54.53. `pg_statistic`

The catalog `pg_statistic` stores statistical data about the contents of the database. Entries are created by `ANALYZE` and subsequently used by the query planner. Note that all the statistical data is inherently approximate, even assuming that it is up-to-date.

Normally there is one entry, with `stainherit = false`, for each table column that has been analyzed. If the table has inheritance children or partitions, a second entry with `stainherit = true` is also created. This row represents the column's statistics over the inheritance tree, i.e., statistics for the data you'd see with `SELECT column FROM table*`, whereas the `stainherit = false` row represents the results of `SELECT column FROM ONLY table`.

`pg_statistic` also stores statistical data about the values of index expressions. These are described as if they were actual data columns; in particular, `starelid` references the index. No entry is made for an ordinary non-expression index column, however, since it would be redundant with the entry for the underlying table column. Currently, entries for index expressions always have `stainherit = false`.

Since different kinds of statistics might be appropriate for different kinds of data, `pg_statistic` is designed not to assume very much about what sort of statistics it stores. Only extremely general statistics (such as nullness) are given dedicated columns in `pg_statistic`. Everything else is stored in “slots”, which are groups of associated columns whose content is identified by a code number in one of the slot's columns.

`pg_statistic` should not be readable by the public, since even statistical information about a table's contents might be considered sensitive. (Example: minimum and maximum values of a salary column might be quite interesting.) `pg_stats` is a publicly readable view on `pg_statistic` that only exposes information about those tables that are readable by the current user.

Table 54.53. `pg_statistic` Columns

Column Type	Description
<code>starelid oid</code> (references <code>pg_class .oid</code>)	The table or index that the described column belongs to
<code>staattnum int2</code> (references <code>pg_attribute .attnum</code>)	The number of the described column
<code>stainherit bool</code>	If true, the stats include values from child tables, not just the values in the specified relation
<code>stanullfrac float4</code>	The fraction of the column's entries that are null
<code>stawidth int4</code>	The average stored width, in bytes, of nonnull entries
<code>stadistinct float4</code>	The number of distinct nonnull data values in the column. A value greater than zero is the actual number of distinct values. A value less than zero is the negative of a multiplier for the number of rows in the table; for example, a column in which about 80% of the values are nonnull and each nonnull value appears about twice on average could be represented by <code>stadistinct = -0.4</code> . A zero value means the number of distinct values is unknown.
<code>stakindN int2</code>	A code number indicating the kind of statistics stored in the <i>N</i> th “slot” of the <code>pg_statistic</code> row.
<code>staopN oid</code> (references <code>pg_operator .oid</code>)	An operator used to derive the statistics stored in the <i>N</i> th “slot”. For example, a histogram slot would show the <code><</code> operator that defines the sort order of the data. Zero if the statistics kind does not require an operator.
<code>stacollN oid</code> (references <code>pg_collation .oid</code>)	The collation used to derive the statistics stored in the <i>N</i> th “slot”. For example, a histogram slot for a collatable column would show the collation that defines the sort order of the data. Zero for noncollatable data.
<code>stanumbersN float4[]</code>	

Column Type	Description
	Numerical statistics of the appropriate kind for the <i>n</i> th “slot”, or null if the slot kind does not involve numerical values
<code>stavalues</code> <i>N</i> <code>anyarray</code>	Column data values of the appropriate kind for the <i>n</i> th “slot”, or null if the slot kind does not store any data values. Each array's element values are actually of the specific column's data type, or a related type such as an array's element type, so there is no way to define these columns' type more specifically than <code>anyarray</code> .

54.54. `pg_statistic_ext`

The catalog `pg_statistic_ext` holds definitions of extended planner statistics. Each row in this catalog corresponds to a *statistics object* created with `CREATE STATISTICS`.

Table 54.54. `pg_statistic_ext` Columns

Column Type	Description
<code>oid</code> <code>oid</code>	Row identifier
<code>stxrelid</code> <code>oid</code> (references <code>pg_class</code> .oid)	Table containing the columns described by this object
<code>stxname</code> <code>name</code>	Name of the statistics object
<code>stxnamespace</code> <code>oid</code> (references <code>pg_namespace</code> .oid)	The OID of the namespace that contains this statistics object
<code>stxowner</code> <code>oid</code> (references <code>pg_authid</code> .oid)	Owner of the statistics object
<code>stxkeys</code> <code>int2</code> <code>vector</code> (references <code>pg_attribute</code> .attnum)	An array of attribute numbers, indicating which table columns are covered by this statistics object; for example a value of <code>1 3</code> would mean that the first and the third table columns are covered
<code>stxstattarget</code> <code>int2</code>	<code>stxstattarget</code> controls the level of detail of statistics accumulated for this statistics object by <code>ANALYZE</code> . A zero value indicates that no statistics should be collected. A null value says to use the maximum of the statistics targets of the referenced columns, if set, or the system default statistics target. Positive values of <code>stxstattarget</code> determine the target number of “most common values” to collect.
<code>stxkind</code> <code>char</code> []	An array containing codes for the enabled statistics kinds; valid values are: <code>d</code> for n-distinct statistics, <code>f</code> for functional dependency statistics, <code>m</code> for most common values (MCV) list statistics, and <code>e</code> for expression statistics
<code>stxexprs</code> <code>pg_node_tree</code>	Expression trees (in <code>nodeToString()</code> representation) for statistics object attributes that are not simple column references. This is a list with one element per expression. Null if all statistics object attributes are simple references.

The `pg_statistic_ext` entry is filled in completely during `CREATE STATISTICS`, but the actual statistical values are not computed then. Subsequent `ANALYZE` commands compute the desired values and populate an entry in the `pg_statistic_ext_data` catalog.

54.55. `pg_statistic_ext_data`

The catalog `pg_statistic_ext_data` holds data for extended planner statistics defined in `pg_statistic_ext`. Each row in this catalog corresponds to a *statistics object* created with `CREATE STATISTICS`.

Normally there is one entry, with `stxdinherit = false`, for each statistics object that has been analyzed. If the table has inheritance children or partitions, a second entry with `stxdinherit = true` is also created. This row represents the statistics object over the inheritance tree, i.e., statistics for the data you'd see with `SELECT * FROM table*`, whereas the `stxdinherit = false` row represents the results of `SELECT * FROM ONLY table`.

Like `pg_statistic`, `pg_statistic_ext_data` should not be readable by the public, since the contents might be considered sensitive. (Example: most common combinations of values in columns might be quite interesting.) `pg_stats_ext` is a publicly readable view on `pg_statistic_ext_data` (after joining with `pg_statistic_ext`) that only exposes information about tables the current user owns.

Table 54.55. `pg_statistic_ext_data` Columns

Column Type	Description
<code>stxoid oid</code> (references <code>pg_statistic_ext .oid</code>)	Extended statistics object containing the definition for this data
<code>stxdinherit bool</code>	If true, the stats include values from child tables, not just the values in the specified relation
<code>stxdndistinct pg_ndistinct</code>	N-distinct counts, serialized as <code>pg_ndistinct type</code>
<code>stxddependencies pg_dependencies</code>	Functional dependency statistics, serialized as <code>pg_dependencies type</code>
<code>stxdmcv pg_mcv_list</code>	MCV (most-common values) list statistics, serialized as <code>pg_mcv_list type</code>
<code>stxdexpr pg_statistic[]</code>	Per-expression statistics, serialized as an array of <code>pg_statistic type</code>

54.56. `pg_subscription`

The catalog `pg_subscription` contains all existing logical replication subscriptions. For more information about logical replication see [Chapter 29](#).

Unlike most system catalogs, `pg_subscription` is shared across all databases of a cluster: there is only one copy of `pg_subscription` per cluster, not one per database.

Access to the column `subconninfo` is revoked from normal users, because it could contain plain-text passwords.

Table 54.56. `pg_subscription` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>subdbid oid</code> (references <code>pg_database .oid</code>)	OID of the database that the subscription resides in
<code>subskiplsn pg_lsn</code>	Finish LSN of the transaction whose changes are to be skipped, if a valid LSN; otherwise 0/0.
<code>subname name</code>	Name of the subscription
<code>subowner oid</code> (references <code>pg_authid .oid</code>)	

Column Type	Description
	Owner of the subscription
subenabled bool	If true, the subscription is enabled and should be replicating
subbinary bool	If true, the subscription will request that the publisher send data in binary format
substream char	Controls how to handle the streaming of in-progress transactions: <code>f</code> = disallow streaming of in-progress transactions, <code>t</code> = spill the changes of in-progress transactions to disk and apply at once after the transaction is committed on the publisher and received by the subscriber, <code>p</code> = apply changes directly using a parallel apply worker if available (same as <code>t</code> if no worker is available)
subtwophasestate char	State codes for two-phase mode: <code>d</code> = disabled, <code>p</code> = pending enablement, <code>e</code> = enabled
subdisableonerr bool	If true, the subscription will be disabled if one of its workers detects an error
subpasswordrequired bool	If true, the subscription will be required to specify a password for authentication
subrunasowner bool	If true, the subscription will be run with the permissions of the subscription owner
subfailover bool	If true, the associated replication slots (i.e. the main slot and the table sync slots) in the upstream database are enabled to be synchronized to the standbys
subconninfo text	Connection string to the upstream database
subslotname name	Name of the replication slot in the upstream database (also used for the local replication origin name); null represents <code>NONE</code>
subsynchronouscommit text	The <code>synchronous_commit</code> setting for the subscription's workers to use
subpublications text[]	Array of subscribed publication names. These reference publications defined in the upstream database. For more on publications see Section 29.1 .
suborigin text	The origin value must be either <code>none</code> or <code>any</code> . The default is <code>any</code> . If <code>none</code> , the subscription will request the publisher to only send changes that don't have an origin. If <code>any</code> , the publisher sends changes regardless of their origin.

54.57. pg_subscription_rel

The catalog `pg_subscription_rel` contains the state for each replicated relation in each subscription. This is a many-to-many mapping.

This catalog only contains tables known to the subscription after running either `CREATE SUBSCRIPTION` or `ALTER SUBSCRIPTION ... REFRESH PUBLICATION`.

Table 54.57. pg_subscription_rel Columns

Column Type	Description
srsubid oid	(references <code>pg_subscription.oid</code>)

Column Type	Description
	Reference to subscription
srrelid oid (references pg_class .oid)	Reference to relation
srsubstate char	State code: <i>i</i> = initialize, <i>d</i> = data is being copied, <i>f</i> = finished table copy, <i>s</i> = synchronized, <i>r</i> = ready (normal replication)
srsublsn pg_lsn	Remote LSN of the state change used for synchronization coordination when in <i>s</i> or <i>r</i> states, otherwise null

54.58. pg_tablespace

The catalog `pg_tablespace` stores information about the available tablespaces. Tables can be placed in particular tablespaces to aid administration of disk layout.

Unlike most system catalogs, `pg_tablespace` is shared across all databases of a cluster: there is only one copy of `pg_tablespace` per cluster, not one per database.

Table 54.58. pg_tablespace Columns

Column Type	Description
oid oid	Row identifier
spcname name	Tablespace name
spcowner oid (references pg_authid .oid)	Owner of the tablespace, usually the user who created it
spcacl aclitem[]	Access privileges; see Section 5.8 for details
spcoptions text[]	Tablespace-level options, as “keyword=value” strings

54.59. pg_transform

The catalog `pg_transform` stores information about transforms, which are a mechanism to adapt data types to procedural languages. See [CREATE TRANSFORM](#) for more information.

Table 54.59. pg_transform Columns

Column Type	Description
oid oid	Row identifier
trftype oid (references pg_type .oid)	OID of the data type this transform is for
trflang oid (references pg_language .oid)	OID of the language this transform is for
trffromsql regproc (references pg_proc .oid)	The OID of the function to use when converting the data type for input to the procedural language (e.g., function parameters). Zero is stored if the default behavior should be used.

Column Type	Description
<code>trftosql regproc</code> (references <code>pg_proc .oid</code>)	The OID of the function to use when converting output from the procedural language (e.g., return values) to the data type. Zero is stored if the default behavior should be used.

54.60. `pg_trigger`

The catalog `pg_trigger` stores triggers on tables and views. See [CREATE TRIGGER](#) for more information.

Table 54.60. `pg_trigger` Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>tgrelid oid</code> (references <code>pg_class .oid</code>)	The table this trigger is on
<code>tgparentid oid</code> (references <code>pg_trigger .oid</code>)	Parent trigger that this trigger is cloned from (this happens when partitions are created or attached to a partitioned table); zero if not a clone
<code>tgname name</code>	Trigger name (must be unique among triggers of same table)
<code>tgfoid oid</code> (references <code>pg_proc .oid</code>)	The function to be called
<code>tgtype int2</code>	Bit mask identifying trigger firing conditions
<code>tgenabled char</code>	Controls in which session_replication_role modes the trigger fires. <code>O</code> = trigger fires in “origin” and “local” modes, <code>D</code> = trigger is disabled, <code>R</code> = trigger fires in “replica” mode, <code>A</code> = trigger fires always.
<code>tgisinternal bool</code>	True if trigger is internally generated (usually, to enforce the constraint identified by <code>tgconstraint</code>)
<code>tgconstrrelid oid</code> (references <code>pg_class .oid</code>)	The table referenced by a referential integrity constraint (zero if trigger is not for a referential integrity constraint)
<code>tgconstrindid oid</code> (references <code>pg_class .oid</code>)	The index supporting a unique, primary key, referential integrity, or exclusion constraint (zero if trigger is not for one of these types of constraint)
<code>tgconstraint oid</code> (references <code>pg_constraint .oid</code>)	The <code>pg_constraint</code> entry associated with the trigger (zero if trigger is not for a constraint)
<code>tgdeferrable bool</code>	True if constraint trigger is deferrable
<code>tginitdeferred bool</code>	True if constraint trigger is initially deferred
<code>tgargs int2</code>	Number of argument strings passed to trigger function
<code>tgattr int2vector</code> (references <code>pg_attribute .attnum</code>)	Column numbers, if trigger is column-specific; otherwise an empty array

Column Type	Description
tgargs bytea	Argument strings to pass to trigger, each NULL-terminated
tgqual pg_node_tree	Expression tree (in <code>nodeToString()</code> representation) for the trigger's <code>WHEN</code> condition, or null if none
tgoldtable name	REFERENCING clause name for <code>OLD TABLE</code> , or null if none
tgnewtable name	REFERENCING clause name for <code>NEW TABLE</code> , or null if none

Currently, column-specific triggering is supported only for `UPDATE` events, and so `tgattr` is relevant only for that event type. `tgtype` might contain bits for other event types as well, but those are presumed to be table-wide regardless of what is in `tgattr`.

Note

When `tgconstraint` is nonzero, `tgconstrrelid`, `tgconstrindid`, `tgdeferrable`, and `tginitdeferred` are largely redundant with the referenced `pg_constraint` entry. However, it is possible for a non-deferrable trigger to be associated with a deferrable constraint: foreign key constraints can have some deferrable and some non-deferrable triggers.

Note

`pg_class.relhastriggers` must be true if a relation has any triggers in this catalog.

54.61. pg_ts_config

The `pg_ts_config` catalog contains entries representing text search configurations. A configuration specifies a particular text search parser and a list of dictionaries to use for each of the parser's output token types. The parser is shown in the `pg_ts_config` entry, but the token-to-dictionary mapping is defined by subsidiary entries in `pg_ts_config_map`.

Postgres Pro's text search features are described at length in [Chapter 12](#).

Table 54.61. `pg_ts_config` Columns

Column Type	Description
oid oid	Row identifier
cfgname name	Text search configuration name
cfgnamespace oid (references <code>pg_namespace</code> .oid)	The OID of the namespace that contains this configuration
cfgowner oid (references <code>pg_authid</code> .oid)	Owner of the configuration
cfgparser oid (references <code>pg_ts_parser</code> .oid)	

Column Type
Description
The OID of the text search parser for this configuration

54.62. pg_ts_config_map

The `pg_ts_config_map` catalog contains entries showing which text search dictionaries should be consulted, and in what order, for each output token type of each text search configuration's parser.

Postgres Pro's text search features are described at length in [Chapter 12](#).

Table 54.62. pg_ts_config_map Columns

Column Type
Description
mapcfg oid (references <code>pg_ts_config</code> .oid) The OID of the <code>pg_ts_config</code> entry owning this map entry
maptokentype int4 A token type emitted by the configuration's parser
mapseqno int4 Order in which to consult this entry (lower <code>mapseqnos</code> first)
mapdict oid (references <code>pg_ts_dict</code> .oid) The OID of the text search dictionary to consult

54.63. pg_ts_dict

The `pg_ts_dict` catalog contains entries defining text search dictionaries. A dictionary depends on a text search template, which specifies all the implementation functions needed; the dictionary itself provides values for the user-settable parameters supported by the template. This division of labor allows dictionaries to be created by unprivileged users. The parameters are specified by a text string `dictinitoption`, whose format and meaning vary depending on the template.

Postgres Pro's text search features are described at length in [Chapter 12](#).

Table 54.63. pg_ts_dict Columns

Column Type
Description
oid oid Row identifier
dictname name Text search dictionary name
dictnamespace oid (references <code>pg_namespace</code> .oid) The OID of the namespace that contains this dictionary
dictowner oid (references <code>pg_authid</code> .oid) Owner of the dictionary
dicttemplate oid (references <code>pg_ts_template</code> .oid) The OID of the text search template for this dictionary
dictinitoption text Initialization option string for the template

54.64. pg_ts_parser

The `pg_ts_parser` catalog contains entries defining text search parsers. A parser is responsible for splitting input text into lexemes and assigning a token type to each lexeme. Since a parser must be implemented by C-language-level functions, creation of new parsers is restricted to database superusers.

Postgres Pro's text search features are described at length in [Chapter 12](#).

Table 54.64. `pg_ts_parser` Columns

Column Type	Description
oid oid	Row identifier
prsname name	Text search parser name
prsnamespace oid (references pg_namespace .oid)	The OID of the namespace that contains this parser
prsstart regproc (references pg_proc .oid)	OID of the parser's startup function
prstoken regproc (references pg_proc .oid)	OID of the parser's next-token function
prsend regproc (references pg_proc .oid)	OID of the parser's shutdown function
prsheadline regproc (references pg_proc .oid)	OID of the parser's headline function (zero if none)
prsllextype regproc (references pg_proc .oid)	OID of the parser's lextype function

54.65. `pg_ts_template`

The `pg_ts_template` catalog contains entries defining text search templates. A template is the implementation skeleton for a class of text search dictionaries. Since a template must be implemented by C-language-level functions, creation of new templates is restricted to database superusers.

Postgres Pro's text search features are described at length in [Chapter 12](#).

Table 54.65. `pg_ts_template` Columns

Column Type	Description
oid oid	Row identifier
tmplname name	Text search template name
tmplnamespace oid (references pg_namespace .oid)	The OID of the namespace that contains this template
tmplinit regproc (references pg_proc .oid)	OID of the template's initialization function (zero if none)
tmpllexize regproc (references pg_proc .oid)	OID of the template's lexize function

54.66. `pg_type`

The catalog `pg_type` stores information about data types. Base types and enum types (scalar types) are created with [CREATE TYPE](#), and domains with [CREATE DOMAIN](#). A composite type is automatically created for each table in the database, to represent the row structure of the table. It is also possible to create composite types with `CREATE TYPE AS`.

Table 54.66. pg_type Columns

Column Type	Description
oid oid	Row identifier
typname name	Data type name
typnamespace oid (references pg_namespace .oid)	The OID of the namespace that contains this type
typowner oid (references pg_authid .oid)	Owner of the type
typlen int2	For a fixed-size type, <code>typlen</code> is the number of bytes in the internal representation of the type. But for a variable-length type, <code>typlen</code> is negative. -1 indicates a “varlena” type (one that has a length word), -2 indicates a null-terminated C string.
typbyval bool	<code>typbyval</code> determines whether internal routines pass a value of this type by value or by reference. <code>typbyval</code> had better be false if <code>typlen</code> is not 1, 2, or 4 (or 8 on machines where Datum is 8 bytes). Variable-length types are always passed by reference. Note that <code>typbyval</code> can be false even if the length would allow pass-by-value.
typtype char	<code>typtype</code> is <code>b</code> for a base type, <code>c</code> for a composite type (e.g., a table's row type), <code>d</code> for a domain, <code>e</code> for an enum type, <code>p</code> for a pseudo-type, <code>r</code> for a range type, or <code>m</code> for a multirange type. See also <code>typrelid</code> and <code>typbasetype</code> .
typcategory char	<code>typcategory</code> is an arbitrary classification of data types that is used by the parser to determine which implicit casts should be “preferred”. See Table 54.67 .
typispreferred bool	True if the type is a preferred cast target within its <code>typcategory</code>
typisdefined bool	True if the type is defined, false if this is a placeholder entry for a not-yet-defined type. When <code>typisdefined</code> is false, nothing except the type name, namespace, and OID can be relied on.
typdelim char	Character that separates two values of this type when parsing array input. Note that the delimiter is associated with the array element data type, not the array data type.
typrelid oid (references pg_class .oid)	If this is a composite type (see <code>typtype</code>), then this column points to the pg_class entry that defines the corresponding table. (For a free-standing composite type, the pg_class entry doesn't really represent a table, but it is needed anyway for the type's pg_attribute entries to link to.) Zero for non-composite types.
typsubscript regproc (references pg_proc .oid)	Subscripting handler function's OID, or zero if this type doesn't support subscripting. Types that are “true” array types have <code>typsubscript</code> = <code>array_subscript_handler</code> , but other types may have other handler functions to implement specialized subscripting behavior.
typelem oid (references pg_type .oid)	If <code>typelem</code> is not zero then it identifies another row in <code>pg_type</code> , defining the type yielded by subscripting. This should be zero if <code>typsubscript</code> is zero. However, it can be zero when <code>typsubscript</code> isn't zero, if the handler doesn't need <code>typelem</code> to determine the subscripting result type. Note that a <code>typelem</code> dependency is considered to imply physical containment of the element type in this type; so DDL changes on the element type might be restricted by the presence of this type.

Column Type	Description
<code>typarray oid</code> (references <code>pg_type .oid</code>)	If <code>typarray</code> is not zero then it identifies another row in <code>pg_type</code> , which is the “true” array type having this type as element
<code>typinput regproc</code> (references <code>pg_proc .oid</code>)	Input conversion function (text format)
<code>typoutput regproc</code> (references <code>pg_proc .oid</code>)	Output conversion function (text format)
<code>typreceive regproc</code> (references <code>pg_proc .oid</code>)	Input conversion function (binary format), or zero if none
<code>typsend regproc</code> (references <code>pg_proc .oid</code>)	Output conversion function (binary format), or zero if none
<code>typmodin regproc</code> (references <code>pg_proc .oid</code>)	Type modifier input function, or zero if type does not support modifiers
<code>typmodout regproc</code> (references <code>pg_proc .oid</code>)	Type modifier output function, or zero to use the standard format
<code>typanalyze regproc</code> (references <code>pg_proc .oid</code>)	Custom <code>ANALYZE</code> function, or zero to use the standard function
<code>typalign char</code>	<p><code>typalign</code> is the alignment required when storing a value of this type. It applies to storage on disk as well as most representations of the value inside Postgres Pro. When multiple values are stored consecutively, such as in the representation of a complete row on disk, padding is inserted before a datum of this type so that it begins on the specified boundary. The alignment reference is the beginning of the first datum in the sequence. Possible values are:</p> <ul style="list-style-type: none"> • <code>c</code> = char alignment, i.e., no alignment needed. • <code>s</code> = short alignment (2 bytes on most machines). • <code>i</code> = int alignment (4 bytes on most machines). • <code>d</code> = double alignment (8 bytes on many machines, but by no means all). • <code>x</code> = xid alignment (8 bytes on all machines).
<code>typstorage char</code>	<p><code>typstorage</code> tells for varlena types (those with <code>typplen = -1</code>) if the type is prepared for toasting and what the default strategy for attributes of this type should be. Possible values are:</p> <ul style="list-style-type: none"> • <code>p</code> (plain): Values must always be stored plain (non-varlena types always use this value). • <code>e</code> (external): Values can be stored in a secondary “TOAST” relation (if relation has one, see <code>pg_class.reltoastrelid</code>). • <code>m</code> (main): Values can be compressed and stored inline. • <code>x</code> (extended): Values can be compressed and/or moved to a secondary relation. <p><code>x</code> is the usual choice for toast-able types. Note that <code>m</code> values can also be moved out to secondary storage, but only as a last resort (<code>e</code> and <code>x</code> values are moved first).</p>
<code>typnotnull bool</code>	<code>typnotnull</code> represents a not-null constraint on a type. Used for domains only.
<code>typbasetype oid</code> (references <code>pg_type .oid</code>)	If this is a domain (see <code>typtype</code>), then <code>typbasetype</code> identifies the type that this one is based on. Zero if this type is not a domain.
<code>tyttypmod int4</code>	

Column Type	Description
	Domains use <code>typmod</code> to record the <code>typmod</code> to be applied to their base type (-1 if base type does not use a <code>typmod</code>). -1 if this type is not a domain.
<code>typndims int4</code>	<code>typndims</code> is the number of array dimensions for a domain over an array (that is, <code>typbasetype</code> is an array type). Zero for types other than domains over array types.
<code>typcollation oid</code> (references <code>pg_collation .oid</code>)	<code>typcollation</code> specifies the collation of the type. If the type does not support collations, this will be zero. A base type that supports collations will have a nonzero value here, typically <code>DEFAULT_COLLATION_OID</code> . A domain over a collatable type can have a collation OID different from its base type's, if one was specified for the domain.
<code>typdefaultbin pg_node_tree</code>	If <code>typdefaultbin</code> is not null, it is the <code>nodeToString()</code> representation of a default expression for the type. This is only used for domains.
<code>typdefault text</code>	<code>typdefault</code> is null if the type has no associated default value. If <code>typdefaultbin</code> is not null, <code>typdefault</code> must contain a human-readable version of the default expression represented by <code>typdefaultbin</code> . If <code>typdefaultbin</code> is null and <code>typdefault</code> is not, then <code>typdefault</code> is the external representation of the type's default value, which can be fed to the type's input converter to produce a constant.
<code>typacl aclitem[]</code>	Access privileges; see Section 5.8 for details

Note

For fixed-width types used in system tables, it is critical that the size and alignment defined in `pg_type` agree with the way that the compiler will lay out the column in a structure representing a table row.

[Table 54.67](#) lists the system-defined values of `typcategory`. Any future additions to this list will also be upper-case ASCII letters. All other ASCII characters are reserved for user-defined categories.

Table 54.67. `typcategory` Codes

Code	Category
A	Array types
B	Boolean types
C	Composite types
D	Date/time types
E	Enum types
G	Geometric types
I	Network address types
N	Numeric types
P	Pseudo-types
R	Range types
S	String types
T	Timespan types
U	User-defined types

Code	Category
V	Bit-string types
X	unknown type
Z	Internal-use types

54.67. pg_user_mapping

The catalog `pg_user_mapping` stores the mappings from local user to remote. Access to this catalog is restricted from normal users, use the view `pg_user_mappings` instead.

Table 54.68. pg_user_mapping Columns

Column Type	Description
<code>oid oid</code>	Row identifier
<code>umuser oid (references <code>pg_authid</code> .oid)</code>	OID of the local role being mapped, or zero if the user mapping is public
<code>umserver oid (references <code>pg_foreign_server</code> .oid)</code>	The OID of the foreign server that contains this mapping
<code>umoptions text[]</code>	User mapping specific options, as “keyword=value” strings

Chapter 55. System Views

In addition to the system catalogs, Postgres Pro provides a number of built-in views. Some system views provide convenient access to some commonly used queries on the system catalogs. Other views provide access to internal server state.

The information schema ([Chapter 38](#)) provides an alternative set of views which overlap the functionality of the system views. Since the information schema is SQL-standard whereas the views described here are Postgres Pro-specific, it's usually better to use the information schema if it provides all the information you need.

[Table 55.1](#) lists the system views described here. More detailed documentation of each view follows below. There are some additional views that provide access to accumulated statistics; they are described in [Table 27.2](#).

55.1. Overview

[Table 55.1](#) lists the system views. More detailed documentation of each catalog follows below. Except where noted, all the views described here are read-only.

Table 55.1. System Views

View Name	Purpose
<code>pg_available_extensions</code>	available extensions
<code>pg_available_extension_versions</code>	available versions of extensions
<code>pg_backend_memory_contexts</code>	backend memory contexts
<code>pg_config</code>	compile-time configuration parameters
<code>pg_cursors</code>	open cursors
<code>pg_file_settings</code>	summary of configuration file contents
<code>pg_group</code>	groups of database users
<code>pg_hba_file_rules</code>	summary of client authentication configuration file contents
<code>pg_ident_file_mappings</code>	summary of client user name mapping configuration file contents
<code>pg_indexes</code>	indexes
<code>pg_locks</code>	locks currently held or awaited
<code>pg_matviews</code>	materialized views
<code>pg_policies</code>	policies
<code>pg_prepared_statements</code>	prepared statements
<code>pg_autoprepared_statements</code>	autoprepared statements
<code>pg_prepared_xacts</code>	prepared transactions
<code>pg_publication_tables</code>	publications and information of their associated tables
<code>pg_replication_origin_status</code>	information about replication origins, including replication progress
<code>pg_replication_slots</code>	replication slot information
<code>pg_roles</code>	database roles
<code>pg_rules</code>	rules
<code>pg_seclabels</code>	security labels
<code>pg_sequences</code>	sequences

View Name	Purpose
pg_settings	parameter settings
pg_shadow	database users
pg_shmem_allocations	shared memory allocations
pg_stats	planner statistics
pg_stats_ext	extended planner statistics
pg_stats_ext_exprs	extended planner statistics for expressions
pg_stats_vacuum_database	statistics about database vacuuming
pg_stats_vacuum_indexes	statistics about index vacuuming
pg_stats_vacuum_tables	statistics about table vacuuming
pg_tables	tables
pg_timezone_abbrevs	time zone abbreviations
pg_timezone_names	time zone names
pg_user	database users
pg_user_mappings	user mappings
pg_views	views
pg_wait_events	wait events

55.2. pg_available_extensions

The `pg_available_extensions` view lists the extensions that are available for installation. See also the [pg_extension](#) catalog, which shows the extensions currently installed.

Table 55.2. `pg_available_extensions` Columns

Column Type	Description
<code>name name</code>	Extension name
<code>default_version text</code>	Name of default version, or <code>NULL</code> if none is specified
<code>installed_version text</code>	Currently installed version of the extension, or <code>NULL</code> if not installed
<code>comment text</code>	Comment string from the extension's control file

The `pg_available_extensions` view is read-only.

55.3. pg_available_extension_versions

The `pg_available_extension_versions` view lists the specific extension versions that are available for installation. See also the [pg_extension](#) catalog, which shows the extensions currently installed.

Table 55.3. `pg_available_extension_versions` Columns

Column Type	Description
<code>name name</code>	Extension name
<code>version text</code>	

Column Type	Description
Version name	
installed bool	True if this version of this extension is currently installed
superuser bool	True if only superusers are allowed to install this extension (but see <code>trusted</code>)
trusted bool	True if the extension can be installed by non-superusers with appropriate privileges
relocatable bool	True if extension can be relocated to another schema
schema name	Name of the schema that the extension must be installed into, or <code>NULL</code> if partially or fully relocatable
requires name[]	Names of prerequisite extensions, or <code>NULL</code> if none
comment text	Comment string from the extension's control file

The `pg_available_extension_versions` view is read-only.

55.4. `pg_backend_memory_contexts`

The view `pg_backend_memory_contexts` displays all the memory contexts of the server process attached to the current session.

`pg_backend_memory_contexts` contains one row for each memory context.

Table 55.4. `pg_backend_memory_contexts` Columns

Column Type	Description
name text	Name of the memory context
ident text	Identification information of the memory context. This field is truncated at 1024 bytes
parent text	Name of the parent of this memory context
level int4	Distance from <code>TopMemoryContext</code> in context tree
total_bytes int8	Total bytes allocated for this memory context
total_nblocks int8	Total number of blocks allocated for this memory context
free_bytes int8	Free space in bytes
free_chunks int8	Total number of free chunks
used_bytes int8	Used space in bytes

By default, the `pg_backend_memory_contexts` view can be read only by superusers or roles with the privileges of the `pg_read_all_stats` role.

55.5. pg_config

The view `pg_config` describes the compile-time configuration parameters of the currently installed version of Postgres Pro. It is intended, for example, to be used by software packages that want to interface to Postgres Pro to facilitate finding the required header files and libraries. It provides the same basic information as the `pg_config` Postgres Pro client application.

By default, the `pg_config` view can be read only by superusers.

Table 55.5. pg_config Columns

Column Type	Description
name text	The parameter name
setting text	The parameter value

55.6. pg_cursors

The `pg_cursors` view lists the cursors that are currently available. Cursors can be defined in several ways:

- via the `DECLARE` statement in SQL
- via the Bind message in the frontend/backend protocol, as described in [Section 56.2.3](#)
- via the Server Programming Interface (SPI), as described in [Section 48.1](#)

The `pg_cursors` view displays cursors created by any of these means. Cursors only exist for the duration of the transaction that defines them, unless they have been declared `WITH HOLD`. Therefore non-holdable cursors are only present in the view until the end of their creating transaction.

Note

Cursors are used internally to implement some of the components of Postgres Pro, such as procedural languages. Therefore, the `pg_cursors` view might include cursors that have not been explicitly created by the user.

Table 55.6. pg_cursors Columns

Column Type	Description
name text	The name of the cursor
statement text	The verbatim query string submitted to declare this cursor
is_holdable bool	true if the cursor is holdable (that is, it can be accessed after the transaction that declared the cursor has committed); false otherwise
is_binary bool	true if the cursor was declared <code>BINARY</code> ; false otherwise
is_scrollable bool	true if the cursor is scrollable (that is, it allows rows to be retrieved in a nonsequential manner); false otherwise
creation_time timestamptz	The time at which the cursor was declared

The `pg_cursors` view is read-only.

55.7. pg_file_settings

The view `pg_file_settings` provides a summary of the contents of the server's configuration file(s). A row appears in this view for each “name = value” entry appearing in the files, with annotations indicating whether the value could be applied successfully. Additional row(s) may appear for problems not linked to a “name = value” entry, such as syntax errors in the files.

This view is helpful for checking whether planned changes in the configuration files will work, or for diagnosing a previous failure. Note that this view reports on the *current* contents of the files, not on what was last applied by the server. (The `pg_settings` view is usually sufficient to determine that.)

By default, the `pg_file_settings` view can be read only by superusers.

Table 55.7. pg_file_settings Columns

Column Type	Description
sourcefile text	Full path name of the configuration file
sourceline int4	Line number within the configuration file where the entry appears
seqno int4	Order in which the entries are processed (1.. <i>n</i>)
name text	Configuration parameter name
setting text	Value to be assigned to the parameter
applied bool	True if the value can be applied successfully
error text	If not null, an error message indicating why this entry could not be applied

If the configuration file contains syntax errors or invalid parameter names, the server will not attempt to apply any settings from it, and therefore all the `applied` fields will read as false. In such a case there will be one or more rows with non-null `error` fields indicating the problem(s). Otherwise, individual settings will be applied if possible. If an individual setting cannot be applied (e.g., invalid value, or the setting cannot be changed after server start) it will have an appropriate message in the `error` field. Another way that an entry might have `applied = false` is that it is overridden by a later entry for the same parameter name; this case is not considered an error so nothing appears in the `error` field.

See [Section 19.1](#) for more information about the various ways to change run-time parameters.

55.8. pg_group

The view `pg_group` exists for backwards compatibility: it emulates a catalog that existed in PostgreSQL before version 8.1. It shows the names and members of all roles that are marked as not `rolcanlogin`, which is an approximation to the set of roles that are being used as groups.

Table 55.8. pg_group Columns

Column Type	Description
groname name (references <code>pg_authid</code> .rolname)	Name of the group
grosysid oid (references <code>pg_authid</code> .oid)	

Column Type	Description
	ID of this group
grolist oid[]	(references pg_authid .oid) An array containing the IDs of the roles in this group

55.9. pg_hba_file_rules

The view `pg_hba_file_rules` provides a summary of the contents of the client authentication configuration file, `pg_hba.conf`. A row appears in this view for each non-empty, non-comment line in the file, with annotations indicating whether the rule could be applied successfully.

This view can be helpful for checking whether planned changes in the authentication configuration file will work, or for diagnosing a previous failure. Note that this view reports on the *current* contents of the file, not on what was last loaded by the server.

By default, the `pg_hba_file_rules` view can be read only by superusers.

Table 55.9. pg_hba_file_rules Columns

Column Type	Description
rule_number int4	Number of this rule, if valid, otherwise NULL. This indicates the order in which each rule is considered until a match is found during authentication.
file_name text	Name of the file containing this rule
line_number int4	Line number of this rule in <code>file_name</code>
type text	Type of connection
database text[]	List of database name(s) to which this rule applies
user_name text[]	List of user and group name(s) to which this rule applies
address text	Host name or IP address, or one of <code>all</code> , <code>samehost</code> , or <code>samenet</code> , or null for local connections
netmask text	IP address mask, or null if not applicable
auth_method text	Authentication method
options text[]	Options specified for authentication method, if any
error text	If not null, an error message indicating why this line could not be processed

Usually, a row reflecting an incorrect entry will have values for only the `line_number` and `error` fields.

See [Chapter 20](#) for more information about client authentication configuration.

55.10. pg_ident_file_mappings

The view `pg_ident_file_mappings` provides a summary of the contents of the client user name mapping configuration file, `pg_ident.conf`. A row appears in this view for each non-empty, non-comment line in the file, with annotations indicating whether the map could be applied successfully.

This view can be helpful for checking whether planned changes in the authentication configuration file will work, or for diagnosing a previous failure. Note that this view reports on the *current* contents of the file, not on what was last loaded by the server.

By default, the `pg_ident_file_mappings` view can be read only by superusers.

Table 55.10. `pg_ident_file_mappings` Columns

Column Type	Description
<code>map_number int4</code>	Number of this map, in priority order, if valid, otherwise <code>NULL</code>
<code>file_name text</code>	Name of the file containing this map
<code>line_number int4</code>	Line number of this map in <code>file_name</code>
<code>map_name text</code>	Name of the map
<code>sys_name text</code>	Detected user name of the client
<code>pg_username text</code>	Requested Postgres Pro user name
<code>error text</code>	If not <code>NULL</code> , an error message indicating why this line could not be processed

Usually, a row reflecting an incorrect entry will have values for only the `line_number` and `error` fields.

See [Chapter 20](#) for more information about client authentication configuration.

55.11. `pg_indexes`

The view `pg_indexes` provides access to useful information about each index in the database.

Table 55.11. `pg_indexes` Columns

Column Type	Description
<code>schemaname name (references pg_namespace .nspname)</code>	Name of schema containing table and index
<code>tablename name (references pg_class .relname)</code>	Name of table the index is for
<code>indexname name (references pg_class .relname)</code>	Name of index
<code>tablespace name (references pg_tablespace .spcname)</code>	Name of tablespace containing index (null if default for database)
<code>indexdef text</code>	Index definition (a reconstructed CREATE INDEX command)

55.12. `pg_locks`

The view `pg_locks` provides access to information about the locks held by active processes within the database server. See [Chapter 13](#) for more discussion of locking.

`pg_locks` contains one row per active lockable object, requested lock mode, and relevant process. Thus, the same lockable object might appear many times, if multiple processes are holding or waiting for locks on it. However, an object that currently has no locks on it will not appear at all.

There are several distinct types of lockable objects: whole relations (e.g., tables), individual pages of relations, individual tuples of relations, transaction IDs (both virtual and permanent IDs), and general database objects (identified by class OID and object OID, in the same way as in [pg_description](#) or [pg_depend](#)). Also, the right to extend a relation is represented as a separate lockable object, as is the right to update `pg_database.datfrozenxid`. Also, “advisory” locks can be taken on numbers that have user-defined meanings.

Table 55.12. `pg_locks` Columns

Column Type	Description
<code>locktype</code> text	Type of the lockable object: <code>relation</code> , <code>extend</code> , <code>frozenid</code> , <code>page</code> , <code>tuple</code> , <code>transactionid</code> , <code>virtualxid</code> , <code>spectoken</code> , <code>object</code> , <code>userlock</code> , <code>advisory</code> , <code>applytransaction</code> , or <code>shardman</code> . (See also Table 27.11 .)
<code>database oid</code> (references <code>pg_database .oid</code>)	OID of the database in which the lock target exists, or zero if the target is a shared object, or null if the target is a transaction ID
<code>relation oid</code> (references <code>pg_class .oid</code>)	OID of the relation targeted by the lock, or null if the target is not a relation or part of a relation
<code>page</code> int4	Page number targeted by the lock within the relation, or null if the target is not a relation page or tuple
<code>tuple</code> int2	Tuple number targeted by the lock within the page, or null if the target is not a tuple
<code>virtualxid</code> text	Virtual ID of the transaction targeted by the lock, or null if the target is not a virtual transaction ID; see Chapter 67
<code>transactionid</code> xid	ID of the transaction targeted by the lock, or null if the target is not a transaction ID; Chapter 67
<code>classid oid</code> (references <code>pg_class .oid</code>)	OID of the system catalog containing the lock target, or null if the target is not a general database object
<code>objid oid</code> (references any OID column)	OID of the lock target within its system catalog, or null if the target is not a general database object
<code>objsubid</code> int2	Column number targeted by the lock (the <code>classid</code> and <code>objid</code> refer to the table itself), or zero if the target is some other general database object, or null if the target is not a general database object
<code>virtualtransaction</code> text	Virtual ID of the transaction that is holding or awaiting this lock
<code>pid</code> int4	Process ID of the server process holding or awaiting this lock, or null if the lock is held by a prepared transaction
<code>mode</code> text	Name of the lock mode held or desired by this process (see Section 13.3.1 and Section 13.2.3)
<code>granted</code> bool	True if lock is held, false if lock is awaited
<code>fastpath</code> bool	

Column Type	Description
	True if lock was taken via fast path, false if taken via main lock table
<code>waitstart timestampz</code>	Time when the server process started waiting for this lock, or null if the lock is held. Note that this can be null for a very short period of time after the wait started even though <code>granted</code> is false.

`granted` is true in a row representing a lock held by the indicated process. False indicates that this process is currently waiting to acquire this lock, which implies that at least one other process is holding or waiting for a conflicting lock mode on the same lockable object. The waiting process will sleep until the other lock is released (or a deadlock situation is detected). A single process can be waiting to acquire at most one lock at a time.

Throughout running a transaction, a server process holds an exclusive lock on the transaction's virtual transaction ID. If a permanent ID is assigned to the transaction (which normally happens only if the transaction changes the state of the database), it also holds an exclusive lock on the transaction's permanent transaction ID until it ends. When a process finds it necessary to wait specifically for another transaction to end, it does so by attempting to acquire share lock on the other transaction's ID (either virtual or permanent ID depending on the situation). That will succeed only when the other transaction terminates and releases its locks.

Although tuples are a lockable type of object, information about row-level locks is stored on disk, not in memory, and therefore row-level locks normally do not appear in this view. If a process is waiting for a row-level lock, it will usually appear in the view as waiting for the permanent transaction ID of the current holder of that row lock.

A speculative insertion lock consists of a transaction ID and a speculative insertion token. The speculative insertion token is displayed in the `objid` column.

Advisory locks can be acquired on keys consisting of either a single `bigint` value or two integer values. A `bigint` key is displayed with its high-order half in the `classid` column, its low-order half in the `objid` column, and `objsubid` equal to 1. The original `bigint` value can be reassembled with the expression `(classid::bigint << 32) | objid::bigint`. Integer keys are displayed with the first key in the `classid` column, the second key in the `objid` column, and `objsubid` equal to 2. The actual meaning of the keys is up to the user. Advisory locks are local to each database, so the `database` column is meaningful for an advisory lock.

Apply transaction locks are used in parallel mode to apply the transaction in logical replication. The remote transaction ID is displayed in the `transactionid` column. The `objsubid` displays the lock subtype which is 0 for the lock used to synchronize the set of changes, and 1 for the lock used to wait for the transaction to finish to ensure commit order.

`pg_locks` provides a global view of all locks in the database cluster, not only those relevant to the current database. Although its `relation` column can be joined against `pg_class.oid` to identify locked relations, this will only work correctly for relations in the current database (those for which the `database` column is either the current database's OID or zero).

The `pid` column can be joined to the `pid` column of the `pg_stat_activity` view to get more information on the session holding or awaiting each lock, for example

```
SELECT * FROM pg_locks pl LEFT JOIN pg_stat_activity psa
    ON pl.pid = psa.pid;
```

Also, if you are using prepared transactions, the `virtualtransaction` column can be joined to the `transaction` column of the `pg_prepared_xacts` view to get more information on prepared transactions that hold locks. (A prepared transaction can never be waiting for a lock, but it continues to hold the locks it acquired while running.) For example:

```
SELECT * FROM pg_locks pl LEFT JOIN pg_prepared_xacts ppx
ON pl.virtualtransaction = '-1/' || ppx.transaction;
```

While it is possible to obtain information about which processes block which other processes by joining `pg_locks` against itself, this is very difficult to get right in detail. Such a query would have to encode knowledge about which lock modes conflict with which others. Worse, the `pg_locks` view does not expose information about which processes are ahead of which others in lock wait queues, nor information about which processes are parallel workers running on behalf of which other client sessions. It is better to use the `pg_blocking_pids()` function (see [Table 9.69](#)) to identify which process(es) a waiting process is blocked behind.

The `pg_locks` view displays data from both the regular lock manager and the predicate lock manager, which are separate systems; in addition, the regular lock manager subdivides its locks into regular and *fast-path* locks. This data is not guaranteed to be entirely consistent. When the view is queried, data on fast-path locks (with `fastpath = true`) is gathered from each backend one at a time, without freezing the state of the entire lock manager, so it is possible for locks to be taken or released while information is gathered. Note, however, that these locks are known not to conflict with any other lock currently in place. After all backends have been queried for fast-path locks, the remainder of the regular lock manager is locked as a unit, and a consistent snapshot of all remaining locks is collected as an atomic action. After unlocking the regular lock manager, the predicate lock manager is similarly locked and all predicate locks are collected as an atomic action. Thus, with the exception of fast-path locks, each lock manager will deliver a consistent set of results, but as we do not lock both lock managers simultaneously, it is possible for locks to be taken or released after we interrogate the regular lock manager and before we interrogate the predicate lock manager.

Locking the regular and/or predicate lock manager could have some impact on database performance if this view is very frequently accessed. The locks are held only for the minimum amount of time necessary to obtain data from the lock managers, but this does not completely eliminate the possibility of a performance impact.

55.13. pg_matviews

The view `pg_matviews` provides access to useful information about each materialized view in the database.

Table 55.13. `pg_matviews` Columns

Column Type	Description
<code>schemaname name</code> (references <code>pg_namespace .nspname</code>)	Name of schema containing materialized view
<code>matviewname name</code> (references <code>pg_class .relname</code>)	Name of materialized view
<code>matviewowner name</code> (references <code>pg_authid .rolname</code>)	Name of materialized view's owner
<code>tablespace name</code> (references <code>pg_tablespace .spcname</code>)	Name of tablespace containing materialized view (null if default for database)
<code>hasindexes bool</code>	True if materialized view has (or recently had) any indexes
<code>ispopulated bool</code>	True if materialized view is currently populated
<code>definition text</code>	Materialized view definition (a reconstructed <code>SELECT</code> query)

55.14. pg_policies

The view `pg_policies` provides access to useful information about each row-level security policy in the database.

Table 55.14. `pg_policies` Columns

Column Type	Description
<code>schemaname name</code> (references <code>pg_namespace .nspname</code>)	Name of schema containing table policy is on
<code>tablename name</code> (references <code>pg_class .relname</code>)	Name of table policy is on
<code>polICYname name</code> (references <code>pg_policy .polname</code>)	Name of policy
<code>permissive text</code>	Is the policy permissive or restrictive?
<code>roles name[]</code>	The roles to which this policy applies
<code>cmd text</code>	The command type to which the policy is applied
<code>qual text</code>	The expression added to the security barrier qualifications for queries that this policy applies to
<code>with_check text</code>	The expression added to the WITH CHECK qualifications for queries that attempt to add rows to this table

55.15. `pg_prepared_statements`

The `pg_prepared_statements` view displays all the prepared statements that are available in the current session. See [PREPARE](#) for more information about prepared statements.

`pg_prepared_statements` contains one row for each prepared statement. Rows are added to the view when a new prepared statement is created and removed when a prepared statement is released (for example, via the [DEALLOCATE](#) command).

Table 55.15. `pg_prepared_statements` Columns

Column Type	Description
<code>name text</code>	The identifier of the prepared statement
<code>statement text</code>	The query string submitted by the client to create this prepared statement. For prepared statements created via SQL, this is the <code>PREPARE</code> statement submitted by the client. For prepared statements created via the frontend/backend protocol, this is the text of the prepared statement itself.
<code>prepare_time timestamptz</code>	The time at which the prepared statement was created
<code>parameter_types regtype[]</code>	The expected parameter types for the prepared statement in the form of an array of <code>regtype</code> . The OID corresponding to an element of this array can be obtained by casting the <code>regtype</code> value to <code>oid</code> .
<code>result_types regtype[]</code>	

Column Type	Description
	The types of the columns returned by the prepared statement in the form of an array of <code>reg-type</code> . The OID corresponding to an element of this array can be obtained by casting the <code>reg-type</code> value to <code>oid</code> . If the prepared statement does not provide a result (e.g., a DML statement), then this field will be null.
<code>from_sql</code> <code>bool</code>	<code>true</code> if the prepared statement was created via the <code>PREPARE SQL</code> command; <code>false</code> if the statement was prepared via the frontend/backend protocol
<code>generic_plans</code> <code>int8</code>	Number of times generic plan was chosen
<code>custom_plans</code> <code>int8</code>	Number of times custom plan was chosen

The `pg_prepared_statements` view is read-only.

55.16. pg_autoprepared_statements

The `pg_autoprepared_statements` view displays all the autoprepared statements that are available in the current session. See [Section 14.5](#) for more information about autoprepared statements.

Table 55.16. pg_autoprepared_statements Columns

Name	Type	Description
<code>statement</code>	<code>text</code>	The query string submitted by the client from which this prepared statement was created. Note that literals in this string are not replaced with prepared statement placeholders.
<code>parameter_types</code>	<code>regtype[]</code>	The expected parameter types for the autoprepared statement in the form of an array of <code>reg-type</code> . The OID corresponding to an element of this array can be obtained by casting the <code>regtype</code> value to <code>oid</code> .
<code>exec_count</code>	<code>int8</code>	Number of times this statement was executed.

The `pg_autoprepared_statements` view is read only.

55.17. pg_prepared_xacts

The view `pg_prepared_xacts` displays information about transactions that are currently prepared for two-phase commit (see [PREPARE TRANSACTION](#) for details).

`pg_prepared_xacts` contains one row per prepared transaction. An entry is removed when the transaction is committed or rolled back.

Table 55.17. pg_prepared_xacts Columns

Column Type	Description
<code>transaction_xid</code>	Numeric transaction identifier of the prepared transaction
<code>gid</code> <code>text</code>	

Column Type	Description
	Global transaction identifier that was assigned to the transaction
prepared timestamp tz	Time at which the transaction was prepared for commit
owner name (references pg_authid .rolname)	Name of the user that executed the transaction
database name (references pg_database .datname)	Name of the database in which the transaction was executed
state3pc text	State of the three-phase commit

When the `pg_prepared_xacts` view is accessed, the internal transaction manager data structures are momentarily locked, and a copy is made for the view to display. This ensures that the view produces a consistent set of results, while not blocking normal operations longer than necessary. Nonetheless there could be some impact on database performance if this view is frequently accessed.

55.18. `pg_publication_tables`

The view `pg_publication_tables` provides information about the mapping between publications and information of tables they contain. Unlike the underlying catalog `pg_publication_rel`, this view expands publications defined as `FOR ALL TABLES` and `FOR TABLES IN SCHEMA`, so for such publications there will be a row for each eligible table.

Table 55.18. `pg_publication_tables` Columns

Column Type	Description
pubname name (references pg_publication .pubname)	Name of publication
schemaname name (references pg_namespace .nspname)	Name of schema containing table
tablename name (references pg_class .relname)	Name of table
attnames name[] (references pg_attribute .attname)	Names of table columns included in the publication. This contains all the columns of the table when the user didn't specify the column list for the table.
rowfilter text	Expression for the table's publication qualifying condition

55.19. `pg_replication_origin_status`

The `pg_replication_origin_status` view contains information about how far replay for a certain origin has progressed. For more on replication origins see [Chapter 51](#).

Table 55.19. `pg_replication_origin_status` Columns

Column Type	Description
local_id oid (references pg_replication_origin .roident)	internal node identifier
external_id text (references pg_replication_origin .roname)	external node identifier
remote_lsn pg_lsn	

Column Type	Description
	The origin node's LSN up to which data has been replicated.
local_lsn pg_lsn	This node's LSN at which <code>remote_lsn</code> has been replicated. Used to flush commit records before persisting data to disk when using asynchronous commits.

55.20. pg_replication_slots

The `pg_replication_slots` view provides a listing of all replication slots that currently exist on the database cluster, along with their current state.

For more on replication slots, see [Section 26.2.6](#) and [Chapter 50](#).

Table 55.20. pg_replication_slots Columns

Column Type	Description
slot_name name	A unique, cluster-wide identifier for the replication slot
plugin name	The base name of the shared object containing the output plugin this logical slot is using, or null for physical slots.
slot_type text	The slot type: <code>physical</code> or <code>logical</code>
datoid oid (references pg_database .oid)	The OID of the database this slot is associated with, or null. Only logical slots have an associated database.
database name (references pg_database .datname)	The name of the database this slot is associated with, or null. Only logical slots have an associated database.
temporary bool	True if this is a temporary replication slot. Temporary slots are not saved to disk and are automatically dropped on error or when the session has finished.
active bool	True if this slot is currently being streamed
active_pid int4	The process ID of the session streaming data for this slot. <code>NULL</code> if inactive.
xmin xid	The oldest transaction that this slot needs the database to retain. <code>VACUUM</code> cannot remove tuples deleted by any later transaction.
catalog_xmin xid	The oldest transaction affecting the system catalogs that this slot needs the database to retain. <code>VACUUM</code> cannot remove catalog tuples deleted by any later transaction.
restart_lsn pg_lsn	The address (LSN) of oldest WAL which still might be required by the consumer of this slot and thus won't be automatically removed during checkpoints unless this LSN gets behind more than max_slot_wal_keep_size from the current LSN. <code>NULL</code> if the LSN of this slot has never been reserved.
confirmed_flush_lsn pg_lsn	The address (LSN) up to which the logical slot's consumer has confirmed receiving data. Data corresponding to the transactions committed before this LSN is not available anymore. <code>NULL</code> for physical slots.

Column	Type	Description
<code>wal_status</code>	<code>text</code>	<p>Availability of WAL files claimed by this slot. Possible values are:</p> <ul style="list-style-type: none"> <code>reserved</code> means that the claimed files are within <code>max_wal_size</code> . <code>extended</code> means that <code>max_wal_size</code> is exceeded but the files are still retained, either by the replication slot or by <code>wal_keep_size</code> . <code>unreserved</code> means that the slot no longer retains the required WAL files and some of them are to be removed at the next checkpoint. This state can return to <code>reserved</code> or <code>extended</code>. <code>lost</code> means that some required WAL files have been removed and this slot is no longer usable. <p>The last two states are seen only when <code>max_slot_wal_keep_size</code> is non-negative. If <code>restart_lsn</code> is NULL, this field is null.</p>
<code>safe_wal_size</code>	<code>int8</code>	<p>The number of bytes that can be written to WAL such that this slot is not in danger of getting in state "lost". It is NULL for lost slots, as well as if <code>max_slot_wal_keep_size</code> is -1.</p>
<code>two_phase</code>	<code>bool</code>	<p>True if the slot is enabled for decoding prepared transactions. Always false for physical slots.</p>
<code>inactive_since</code>	<code>timestampz</code>	<p>The time when the slot became inactive. NULL if the slot is currently being streamed. Note that for slots on the standby that are being synced from a primary server (whose <code>synced</code> field is <code>true</code>), the <code>inactive_since</code> indicates the time when slot synchronization (see Section 50.2.3) was most recently stopped. NULL if the slot has always been synchronized. On standby, this is useful for slots that are being synced from a primary server (whose <code>synced</code> field is <code>true</code>) so they know when the slot stopped being synchronized.</p>
<code>conflicting</code>	<code>bool</code>	<p>True if this logical slot conflicted with recovery (and so is now invalidated). When this column is true, check <code>invalidation_reason</code> column for the conflict reason. Always NULL for physical slots.</p>
<code>invalidation_reason</code>	<code>text</code>	<p>The reason for the slot's invalidation. It is set for both logical and physical slots. NULL if the slot is not invalidated. Possible values are:</p> <ul style="list-style-type: none"> <code>wal_removed</code> means that the required WAL has been removed. <code>rows_removed</code> means that the required rows have been removed. It is set only for logical slots. <code>wal_level_insufficient</code> means that the primary doesn't have a <code>wal_level</code> sufficient to perform logical decoding. It is set only for logical slots.
<code>failover</code>	<code>bool</code>	<p>True if this is a logical slot enabled to be synced to the standbys so that logical replication can be resumed from the new primary after failover. Always false for physical slots.</p>
<code>synced</code>	<code>bool</code>	<p>True if this is a logical slot that was synced from a primary server. On a hot standby, the slots with the <code>synced</code> column marked as true can neither be used for logical decoding nor dropped manually. The value of this column has no meaning on the primary server; the column value on the primary is default false for all slots but may (if leftover from a promoted standby) also be true.</p>

55.21. pg_roles

The view `pg_roles` provides access to information about database roles. This is simply a publicly readable view of `pg_authid` that blanks out the password field.

Table 55.21. pg_roles Columns

Column	Type	Description
rolname	name	Role name
rolsuper	bool	Role has superuser privileges
rolinherit	bool	Role automatically inherits privileges of roles it is a member of
rolcreaterole	bool	Role can create more roles
rolcreatedb	bool	Role can create databases
rolcanlogin	bool	Role can log in. That is, this role can be given as the initial session authorization identifier
rolreplication	bool	Role is a replication role. A replication role can initiate replication connections and create and drop replication slots.
rolconnlimit	int4	For roles that can log in, this sets maximum number of concurrent connections this role can make. -1 means no limit.
rolprofile	regprofile (references pg_profile .oid)	The OID of the role's profile
rolloginattempts	int4	Number of consecutive failed login attempts of a user. It is always 0 if <code>FAILED_LOGIN_ATTEMPTS</code> value of the corresponding role profile is <code>UNLIMITED</code> . After a successful login, <code>rolloginattempts</code> is reset to 0.
rollastlogin	timestampz	The timestamp the role logged in last time
rolfirstfailedauth	timestampz	The timestamp of the role's first authentication failure
rolstatus	int2	Status of the role: 0 for the active role, 1 the role is manually locked (see <code>ACCOUNT LOCK</code> in ALTER ROLE), 2 the role is locked because of inactivity (see parameter <code>USER_INACTIVE_TIME</code> in CREATE PROFILE), 4 the role is locked because the number of consecutive authentication failures has reached the limit (see parameter <code>FAILED_LOGIN_ATTEMPTS</code> in CREATE PROFILE)
rolpassword	text	Not the password (always reads as <code>*****</code>)
rolvaliduntil	timestampz	Password expiry time (only used for password authentication); null if no expiration
rolpasssetat	timestampz	Password set time (only used for password authentication); null if password is not set.
rolbypassrsls	bool	Role bypasses every row-level security policy, see Section 5.9 for more information.
rolconfig	text[]	Role-specific defaults for run-time configuration variables
oid	oid (references pg_authid .oid)	ID of role

55.22. pg_rules

The view `pg_rules` provides access to useful information about query rewrite rules.

Table 55.22. pg_rules Columns

Column Type	Description
<code>schemaname name</code> (references <code>pg_namespace .nspname</code>)	Name of schema containing table
<code>tablename name</code> (references <code>pg_class .relname</code>)	Name of table the rule is for
<code>rulename name</code> (references <code>pg_rewrite .rulename</code>)	Name of rule
<code>definition text</code>	Rule definition (a reconstructed creation command)

The `pg_rules` view excludes the `ON SELECT` rules of views and materialized views; those can be seen in `pg_views` and `pg_matviews`.

55.23. pg_seclabels

The view `pg_seclabels` provides information about security labels. It is an easier-to-query version of the `pg_seclabel` catalog.

Table 55.23. pg_seclabels Columns

Column Type	Description
<code>objoid oid</code> (references any OID column)	The OID of the object this security label pertains to
<code>classoid oid</code> (references <code>pg_class .oid</code>)	The OID of the system catalog this object appears in
<code>objsubid int4</code>	For a security label on a table column, this is the column number (the <code>objoid</code> and <code>classoid</code> refer to the table itself). For all other object types, this column is zero.
<code>objtype text</code>	The type of object to which this label applies, as text.
<code>objnamespace oid</code> (references <code>pg_namespace .oid</code>)	The OID of the namespace for this object, if applicable; otherwise NULL.
<code>objname text</code>	The name of the object to which this label applies, as text.
<code>provider text</code> (references <code>pg_seclabel .provider</code>)	The label provider associated with this label.
<code>label text</code> (references <code>pg_seclabel .label</code>)	The security label applied to this object.

55.24. pg_sequences

The view `pg_sequences` provides access to useful information about each sequence in the database.

Table 55.24. pg_sequences Columns

Column Type	Description
schemaname name (references pg_namespace .nspname)	Name of schema containing sequence
sequencename name (references pg_class .relname)	Name of sequence
sequenceowner name (references pg_authid .rolname)	Name of sequence's owner
data_type regtype (references pg_type .oid)	Data type of the sequence
start_value int8	Start value of the sequence
min_value int8	Minimum value of the sequence
max_value int8	Maximum value of the sequence
increment_by int8	Increment value of the sequence
cycle bool	Whether the sequence cycles
cache_size int8	Cache size of the sequence
last_value int8	The last sequence value written to disk. If caching is used, this value can be greater than the last value handed out from the sequence.

The `last_value` column will read as null if any of the following are true:

- The sequence has not been read from yet.
- The current user does not have `USAGE` or `SELECT` privilege on the sequence.
- The sequence is unlogged and the server is a standby.

55.25. pg_settings

The view `pg_settings` provides access to run-time parameters of the server. It is essentially an alternative interface to the [SHOW](#) and [SET](#) commands. It also provides access to some facts about each parameter that are not directly available from [SHOW](#), such as minimum and maximum values.

Table 55.25. pg_settings Columns

Column Type	Description
name text	Run-time configuration parameter name
setting text	Current value of the parameter
unit text	Implicit unit of the parameter
category text	Logical group of the parameter
short_desc text	

Column	Type	Description
		A brief description of the parameter
extra_desc	text	Additional, more detailed, description of the parameter
context	text	Context required to set the parameter's value (see below)
vartype	text	Parameter type (bool, enum, integer, real, or string)
source	text	Source of the current parameter value
min_val	text	Minimum allowed value of the parameter (null for non-numeric values)
max_val	text	Maximum allowed value of the parameter (null for non-numeric values)
enumvals	text[]	Allowed values of an enum parameter (null for non-enum values)
boot_val	text	Parameter value assumed at server startup if the parameter is not otherwise set
reset_val	text	Value that RESET would reset the parameter to in the current session
sourcefile	text	Configuration file the current value was set in (null for values set from sources other than configuration files, or when examined by a user who neither is a superuser nor has privileges of <code>pg_read_all_settings</code>); helpful when using <code>include</code> directives in configuration files
sourceline	int4	Line number within the configuration file the current value was set at (null for values set from sources other than configuration files, or when examined by a user who neither is a superuser nor has privileges of <code>pg_read_all_settings</code>).
pending_restart	bool	<code>true</code> if the value has been changed in the configuration file but needs a restart; or <code>false</code> otherwise.

There are several possible values of `context`. In order of decreasing difficulty of changing the setting, they are:

internal

These settings cannot be changed directly; they reflect internally determined values. Some of them may be adjustable by rebuilding the server with different configuration options, or by changing options supplied to `initdb`.

postmaster

These settings can only be applied when the server starts, so any change requires restarting the server. Values for these settings are typically stored in the `postgresql.conf` file, or passed on the command line when starting the server. Of course, settings with any of the lower `context` types can also be set at server start time.

sighup

Changes to these settings can be made in `postgresql.conf` without restarting the server. Send a **SIGHUP** signal to the postmaster to cause it to re-read `postgresql.conf` and apply the changes. The postmaster will also forward the **SIGHUP** signal to its child processes so that they all pick up the new value.

superuser-backend

Changes to these settings can be made in `postgresql.conf` without restarting the server. They can also be set for a particular session in the connection request packet (for example, via libpq's `PGOPTIONS` environment variable), but only if the connecting user is a superuser or has been granted the appropriate `SET` privilege. However, these settings never change in a session after it is started. If you change them in `postgresql.conf`, send a `SIGHUP` signal to the postmaster to cause it to re-read `postgresql.conf`. The new values will only affect subsequently-launched sessions.

backend

Changes to these settings can be made in `postgresql.conf` without restarting the server. They can also be set for a particular session in the connection request packet (for example, via libpq's `PGOPTIONS` environment variable); any user can make such a change for their session. However, these settings never change in a session after it is started. If you change them in `postgresql.conf`, send a `SIGHUP` signal to the postmaster to cause it to re-read `postgresql.conf`. The new values will only affect subsequently-launched sessions.

superuser

These settings can be set from `postgresql.conf`, or within a session via the `SET` command; but only superusers and users with the appropriate `SET` privilege can change them via `SET`. Changes in `postgresql.conf` will affect existing sessions only if no session-local value has been established with `SET`.

user

These settings can be set from `postgresql.conf`, or within a session via the `SET` command. Any user is allowed to change their session-local value. Changes in `postgresql.conf` will affect existing sessions only if no session-local value has been established with `SET`.

See [Section 19.1](#) for more information about the various ways to change these parameters.

This view cannot be inserted into or deleted from, but it can be updated. An `UPDATE` applied to a row of `pg_settings` is equivalent to executing the `SET` command on that named parameter. The change only affects the value used by the current session. If an `UPDATE` is issued within a transaction that is later aborted, the effects of the `UPDATE` command disappear when the transaction is rolled back. Once the surrounding transaction is committed, the effects will persist until the end of the session, unless overridden by another `UPDATE` or `SET`.

This view does not display [customized options](#) unless the extension module that defines them has been loaded by the backend process executing the query (e.g., via a mention in [shared_preload_libraries](#), a call to a C function in the extension, or the `LOAD` command). For example, since [archive modules](#) are normally loaded only by the archiver process not regular sessions, this view will not display any customized options defined by such modules unless special action is taken to load them into the backend process executing the query.

55.26. pg_shadow

The view `pg_shadow` exists for backwards compatibility: it emulates a catalog that existed in PostgreSQL before version 8.1. It shows properties of all roles that are marked as `rolcanlogin` in `pg_authid`.

The name stems from the fact that this table should not be readable by the public since it contains passwords. `pg_user` is a publicly readable view on `pg_shadow` that blanks out the password field.

Table 55.26. pg_shadow Columns

Column Type	Description
username name	(references <code>pg_authid</code> .rolname)

Column Type	Description
User name	
usesysid oid (references <code>pg_authid .oid</code>)	ID of this user
usecreatedb bool	User can create databases
usesuper bool	User is a superuser
userepl bool	User can initiate streaming replication and put the system in and out of backup mode.
usebypassrls bool	User bypasses every row-level security policy, see Section 5.9 for more information.
passwd text	Password (possibly encrypted); null if none. See <code>pg_authid</code> for details of how encrypted passwords are stored.
valuntil timestamptz	Password expiry time (only used for password authentication)
useconfig text[]	Session defaults for run-time configuration variables

55.27. pg_shmem_allocations

The `pg_shmem_allocations` view shows allocations made from the server's main shared memory segment. This includes both memory allocated by Postgres Pro itself and memory allocated by extensions using the mechanisms detailed in [Section 39.10.10](#).

Note that this view does not include memory allocated using the dynamic shared memory infrastructure.

Table 55.27. `pg_shmem_allocations` Columns

Column Type	Description
name text	The name of the shared memory allocation. NULL for unused memory and <code><anonymous></code> for anonymous allocations.
off int8	The offset at which the allocation starts. NULL for anonymous allocations, since details related to them are not known.
size int8	Size of the allocation in bytes
allocated_size int8	Size of the allocation in bytes including padding. For anonymous allocations, no information about padding is available, so the <code>size</code> and <code>allocated_size</code> columns will always be equal. Padding is not meaningful for free memory, so the columns will be equal in that case also.

Anonymous allocations are allocations that have been made with `ShmemAlloc()` directly, rather than via `ShmemInitStruct()` or `ShmemInitHash()`.

By default, the `pg_shmem_allocations` view can be read only by superusers or roles with privileges of the `pg_read_all_stats` role.

55.28. pg_stats

The view `pg_stats` provides access to the information stored in the `pg_statistic` catalog. This view allows access only to rows of `pg_statistic` that correspond to tables the user has permission to read, and therefore it is safe to allow public read access to this view.

`pg_stats` is also designed to present the information in a more readable format than the underlying catalog — at the cost that its schema must be extended whenever new slot types are defined for `pg_statistic`.

Table 55.28. `pg_stats` Columns

Column Type	Description
<code>schemaname name</code> (references <code>pg_namespace .nspname</code>)	Name of schema containing table
<code>tablename name</code> (references <code>pg_class .relname</code>)	Name of table
<code>attname name</code> (references <code>pg_attribute .attname</code>)	Name of column described by this row
<code>inherited bool</code>	If true, this row includes values from child tables, not just the values in the specified table
<code>null_frac float4</code>	Fraction of column entries that are null
<code>avg_width int4</code>	Average width in bytes of column's entries
<code>n_distinct float4</code>	If greater than zero, the estimated number of distinct values in the column. If less than zero, the negative of the number of distinct values divided by the number of rows. (The negated form is used when <code>ANALYZE</code> believes that the number of distinct values is likely to increase as the table grows; the positive form is used when the column seems to have a fixed number of possible values.) For example, -1 indicates a unique column in which the number of distinct values is the same as the number of rows.
<code>most_common_vals anyarray</code>	A list of the most common values in the column. (Null if no values seem to be more common than any others.)
<code>most_common_freqs float4[]</code>	A list of the frequencies of the most common values, i.e., number of occurrences of each divided by total number of rows. (Null when <code>most_common_vals</code> is.)
<code>histogram_bounds anyarray</code>	A list of values that divide the column's values into groups of approximately equal population. The values in <code>most_common_vals</code> , if present, are omitted from this histogram calculation. (This column is null if the column data type does not have a < operator or if the <code>most_common_vals</code> list accounts for the entire population.)
<code>correlation float4</code>	Statistical correlation between physical row ordering and logical ordering of the column values. This ranges from -1 to +1. When the value is near -1 or +1, an index scan on the column will be estimated to be cheaper than when it is near zero, due to reduction of random access to the disk. (This column is null if the column data type does not have a < operator.)
<code>most_common_elems anyarray</code>	A list of non-null element values most often appearing within values of the column. (Null for scalar types.)
<code>most_common_elem_freqs float4[]</code>	A list of the frequencies of the most common element values, i.e., the fraction of rows containing at least one instance of the given value. Two or three additional values follow the per-el-

Column Type	Description
	ement frequencies; these are the minimum and maximum of the preceding per-element frequencies, and optionally the frequency of null elements. (Null when <code>most_common_elems</code> is.)
<code>elem_count_histogram float4[]</code>	A histogram of the counts of distinct non-null element values within the values of the column, followed by the average number of distinct non-null elements. (Null for scalar types.)
<code>range_length_histogram anyarray</code>	A histogram of the lengths of non-empty and non-null range values of a range type column. (Null for non-range types.) This histogram is calculated using the <code>subtype_diff</code> range function regardless of whether range bounds are inclusive.
<code>range_empty_frac float4</code>	Fraction of column entries whose values are empty ranges. (Null for non-range types.)
<code>range_bounds_histogram anyarray</code>	A histogram of lower and upper bounds of non-empty and non-null range values. (Null for non-range types.) These two histograms are represented as a single array of ranges, whose lower bounds represent the histogram of lower bounds, and upper bounds represent the histogram of upper bounds.

The maximum number of entries in the array fields can be controlled on a column-by-column basis using the `ALTER TABLE SET STATISTICS` command, or globally by setting the `default_statistics_target` run-time parameter.

55.29. pg_stats_ext

The view `pg_stats_ext` provides access to information about each extended statistics object in the database, combining information stored in the `pg_statistic_ext` and `pg_statistic_ext_data` catalogs. This view allows access only to rows of `pg_statistic_ext` and `pg_statistic_ext_data` that correspond to tables the user owns, and therefore it is safe to allow public read access to this view.

`pg_stats_ext` is also designed to present the information in a more readable format than the underlying catalogs — at the cost that its schema must be extended whenever new types of extended statistics are added to `pg_statistic_ext`.

Table 55.29. pg_stats_ext Columns

Column Type	Description
<code>schemaname name (references pg_namespace .nspname)</code>	Name of schema containing table
<code>tablename name (references pg_class .relname)</code>	Name of table
<code>statistics_schemaname name (references pg_namespace .nspname)</code>	Name of schema containing extended statistics object
<code>statistics_name name (references pg_statistic_ext .stxname)</code>	Name of extended statistics object
<code>statistics_owner name (references pg_authid .rolname)</code>	Owner of the extended statistics object
<code>attnames name[] (references pg_attribute .attname)</code>	Names of the columns included in the extended statistics object
<code>exprs text []</code>	

Column Type	Description
	Expressions included in the extended statistics object
kinds char[]	Types of extended statistics object enabled for this record
inherited bool (references <code>pg_statistic_ext_data</code> . <code>stxdinherit</code>)	If true, the stats include values from child tables, not just the values in the specified relation
n_distinct pg_ndistinct	N-distinct counts for combinations of column values. If greater than zero, the estimated number of distinct values in the combination. If less than zero, the negative of the number of distinct values divided by the number of rows. (The negated form is used when <code>ANALYZE</code> believes that the number of distinct values is likely to increase as the table grows; the positive form is used when the column seems to have a fixed number of possible values.) For example, -1 indicates a unique combination of columns in which the number of distinct combinations is the same as the number of rows.
dependencies pg_dependencies	Functional dependency statistics
most_common_vals text[]	A list of the most common combinations of values in the columns. (Null if no combinations seem to be more common than any others.)
most_common_val_nulls bool[]	A list of NULL flags for the most common combinations of values. (Null when <code>most_common_vals</code> is.)
most_common_freqs float8[]	A list of the frequencies of the most common combinations, i.e., number of occurrences of each divided by total number of rows. (Null when <code>most_common_vals</code> is.)
most_common_base_freqs float8[]	A list of the base frequencies of the most common combinations, i.e., product of per-value frequencies. (Null when <code>most_common_vals</code> is.)

The maximum number of entries in the array fields can be controlled on a column-by-column basis using the `ALTER TABLE SET STATISTICS` command, or globally by setting the `default_statistics_target` run-time parameter.

55.30. pg_stats_ext_exprs

The view `pg_stats_ext_exprs` provides access to information about all expressions included in extended statistics objects, combining information stored in the `pg_statistic_ext` and `pg_statistic_ext_data` catalogs. This view allows access only to rows of `pg_statistic_ext` and `pg_statistic_ext_data` that correspond to tables the user owns, and therefore it is safe to allow public read access to this view.

`pg_stats_ext_exprs` is also designed to present the information in a more readable format than the underlying catalogs — at the cost that its schema must be extended whenever the structure of statistics in `pg_statistic_ext` changes.

Table 55.30. pg_stats_ext_exprs Columns

Column Type	Description
schemaname name (references <code>pg_namespace</code> . <code>nspname</code>)	Name of schema containing table
tablename name (references <code>pg_class</code> . <code>relname</code>)	Name of table the statistics object is defined on
statistics_schemaname name (references <code>pg_namespace</code> . <code>nspname</code>)	

Column Type	Description
	Name of schema containing extended statistics object
statistics_name name (references pg_statistic_ext .stxname)	Name of extended statistics object
statistics_owner name (references pg_authid .rolname)	Owner of the extended statistics object
expr text	Expression included in the extended statistics object
inherited bool (references pg_statistic_ext_data .stxdinherit)	If true, the stats include values from child tables, not just the values in the specified relation
null_frac float4	Fraction of expression entries that are null
avg_width int4	Average width in bytes of expression's entries
n_distinct float4	If greater than zero, the estimated number of distinct values in the expression. If less than zero, the negative of the number of distinct values divided by the number of rows. (The negated form is used when <code>ANALYZE</code> believes that the number of distinct values is likely to increase as the table grows; the positive form is used when the expression seems to have a fixed number of possible values.) For example, -1 indicates a unique expression in which the number of distinct values is the same as the number of rows.
most_common_vals anyarray	A list of the most common values in the expression. (Null if no values seem to be more common than any others.)
most_common_freqs float4[]	A list of the frequencies of the most common values, i.e., number of occurrences of each divided by total number of rows. (Null when <code>most_common_vals</code> is.)
histogram_bounds anyarray	A list of values that divide the expression's values into groups of approximately equal population. The values in <code>most_common_vals</code> , if present, are omitted from this histogram calculation. (This expression is null if the expression data type does not have a <code><</code> operator or if the <code>most_common_vals</code> list accounts for the entire population.)
correlation float4	Statistical correlation between physical row ordering and logical ordering of the expression values. This ranges from -1 to +1. When the value is near -1 or +1, an index scan on the expression will be estimated to be cheaper than when it is near zero, due to reduction of random access to the disk. (This expression is null if the expression's data type does not have a <code><</code> operator.)
most_common_elems anyarray	A list of non-null element values most often appearing within values of the expression. (Null for scalar types.)
most_common_elem_freqs float4[]	A list of the frequencies of the most common element values, i.e., the fraction of rows containing at least one instance of the given value. Two or three additional values follow the per-element frequencies; these are the minimum and maximum of the preceding per-element frequencies, and optionally the frequency of null elements. (Null when <code>most_common_elems</code> is.)
elem_count_histogram float4[]	A histogram of the counts of distinct non-null element values within the values of the expression, followed by the average number of distinct non-null elements. (Null for scalar types.)

The maximum number of entries in the array fields can be controlled on a column-by-column basis using the `ALTER TABLE SET STATISTICS` command, or globally by setting the `default_statistics_target` run-time parameter.

55.31. pg_stats_vacuum_database

The view `pg_stats_vacuum_database` will contain one row for each database in the current cluster, showing statistics about vacuuming that database.

Table 55.31. pg_stats_vacuum_database Columns

Column	Type	Description
<code>dboid</code>	<code>oid</code>	OID of a database
<code>db_blks_read</code>	<code>int8</code>	Number of database blocks read by vacuum operations performed on this database
<code>db_blks_hit</code>	<code>int8</code>	Number of times database blocks were found in the buffer cache by vacuum operations performed on this database
<code>total_blks_dirtied</code>	<code>int8</code>	Number of database blocks dirtied by vacuum operations performed on this database
<code>total_blks_written</code>	<code>int8</code>	Number of database blocks written by vacuum operations performed on this database
<code>wal_records</code>	<code>int8</code>	Total number of WAL records generated by vacuum operations performed on this database
<code>wal_fpi</code>	<code>int8</code>	Total number of WAL full page images generated by vacuum operations performed on this database
<code>wal_bytes</code>	<code>numeric</code>	Total amount of WAL bytes generated by vacuum operations performed on this database
<code>blk_read_time</code>	<code>float8</code>	Time spent reading database blocks by vacuum operations performed on this database, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>blk_write_time</code>	<code>float8</code>	Time spent writing database blocks by vacuum operations performed on this database, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>delay_time</code>	<code>float8</code>	Time spent sleeping in a vacuum delay point by vacuum operations performed on this database, in milliseconds (see Section 19.4.4 for details)
<code>system_time</code>	<code>float8</code>	System CPU time of vacuuming this database, in milliseconds
<code>user_time</code>	<code>float8</code>	User CPU time of vacuuming this database, in milliseconds
<code>total_time</code>	<code>float8</code>	Total time of vacuuming this database, in milliseconds
<code>interrupts</code>	<code>int4</code>	Number of times vacuum operations performed on this database were interrupted on any errors

Note

On upgrading your server just by replacing the executables, update the `pg_stats_vacuum_database` view as superuser who ran `initdb`:

- For each database, including `template1`, but not including `template0`, run:

```
CREATE OR REPLACE VIEW pg_stats_vacuum_database AS
SELECT
    db.oid as dboid,  stats.db_blks_read,
    stats.db_blks_hit,
    stats.total_blks_dirtied,
    stats.total_blks_written,  stats.wal_records,
    stats.wal_fpi,
    stats.wal_bytes,  stats.blk_read_time,
    stats.blk_write_time,  stats.delay_time,
    stats.system_time,
    stats.user_time,
    stats.total_time,  stats.interrupts
FROM
    pg_database db LEFT JOIN pg_stats_vacuum_database(db.oid) stats
ON
    db.oid = stats.dboid;
```

- For the `template0` database, run:

```
\c template1
ALTER DATABASE template0 ALLOW_CONNECTIONS on;
\c template0
CREATE OR REPLACE VIEW pg_stats_vacuum_database AS
SELECT
    db.oid as dboid,  stats.db_blks_read,
    stats.db_blks_hit,
    stats.total_blks_dirtied,
    stats.total_blks_written,  stats.wal_records,
    stats.wal_fpi,
    stats.wal_bytes,  stats.blk_read_time,
    stats.blk_write_time,  stats.delay_time,
    stats.system_time,
    stats.user_time,
    stats.total_time,  stats.interrupts
FROM
    pg_database db LEFT JOIN pg_stats_vacuum_database(db.oid) stats
ON
    db.oid = stats.dboid;
\c template1
ALTER DATABASE template0 ALLOW_CONNECTIONS off;
```

55.32. `pg_stats_vacuum_indexes`

The view `pg_stats_vacuum_indexes` will contain one row for each index in the current database (including TOAST table indexes), showing statistics about vacuuming that specific index.

Table 55.32. `pg_stats_vacuum_indexes` Columns

Column Type	Description
<code>relid oid</code>	OID of an index

Column	Type	Description
schema_name	name	Name of the schema this index is in
relname	name	Name of this index
total_blks_read	int8	Number of database blocks read by vacuum operations performed on this index
total_blks_hit	int8	Number of times database blocks were found in the buffer cache by vacuum operations performed on this index
total_blks_dirtied	int8	Number of database blocks dirtied by vacuum operations performed on this index
total_blks_written	int8	Number of database blocks written by vacuum operations performed on this index
rel_blks_read	int8	Number of blocks vacuum operations read from this index
rel_blks_hit	int8	Number of times blocks of this index were already found in the buffer cache by vacuum operations, so that a read was not necessary (this only includes hits in the Postgres Pro buffer cache, not the operating system's file system cache)
pages_deleted	int8	Number of pages deleted by vacuum operations performed on this index
tuples_deleted	int8	Number of dead tuples vacuum operations deleted from this index
wal_records	int8	Total number of WAL records generated by vacuum operations performed on this index
wal_fpi	int8	Total number of WAL full page images generated by vacuum operations performed on this index
wal_bytes	numeric	Total amount of WAL bytes generated by vacuum operations performed on this index
blk_read_time	int8	Time spent reading database blocks by vacuum operations performed on this index, in milliseconds (if track_io_timing is enabled, otherwise zero)
blk_write_time	int8	Time spent writing database blocks by vacuum operations performed on this index, in milliseconds (if track_io_timing is enabled, otherwise zero)
delay_time	float8	Time spent sleeping in a vacuum delay point by vacuum operations performed on this index, in milliseconds (see Section 19.4.4 for details)
system_time	float8	System CPU time of vacuuming this index, in milliseconds
user_time	float8	User CPU time of vacuuming this index, in milliseconds
total_time	float8	Total time of vacuuming this index, in milliseconds
interrupts	integer	Number of times vacuum operations performed on this index were interrupted on any errors

55.33. pg_stats_vacuum_tables

The view `pg_stats_vacuum_tables` will contain one row for each table in the current database (including TOAST tables), showing statistics about vacuuming that specific table.

Table 55.33. pg_stats_vacuum_tables Columns

Column	Type	Description
<code>relid</code>	<code>oid</code>	OID of a table
<code>schema</code>	<code>name</code>	Name of the schema this table is in
<code>relname</code>	<code>name</code>	Name of this table
<code>total_blks_read</code>	<code>int8</code>	Number of database blocks read by vacuum operations performed on this table
<code>total_blks_hit</code>	<code>int8</code>	Number of times database blocks were found in the buffer cache by vacuum operations performed on this table
<code>total_blks_dirtied</code>	<code>int8</code>	Number of database blocks dirtied by vacuum operations performed on this table
<code>total_blks_written</code>	<code>int8</code>	Number of database blocks written by vacuum operations performed on this table
<code>rel_blks_read</code>	<code>int8</code>	Number of blocks vacuum operations read from this table
<code>rel_blks_hit</code>	<code>int8</code>	Number of times blocks of this table were already found in the buffer cache by vacuum operations, so that a read was not necessary (this only includes hits in the Postgres Pro buffer cache, not the operating system's file system cache)
<code>pages_scanned</code>	<code>int8</code>	Number of pages examined by vacuum operations performed on this table
<code>pages_removed</code>	<code>int8</code>	Number of pages removed from the physical storage by vacuum operations performed on this table
<code>pages_frozen</code>	<code>int8</code>	Number of times vacuum operations marked pages of this table as all-frozen in the visibility map
<code>pages_all_visible</code>	<code>int8</code>	Number of times vacuum operations marked pages of this table as all-visible in the visibility map
<code>tuples_deleted</code>	<code>int8</code>	Number of dead tuples vacuum operations deleted from this table
<code>tuples_frozen</code>	<code>int8</code>	Number of tuples of this table that vacuum operations marked as frozen
<code>dead_tuples</code>	<code>int8</code>	Number of dead tuples vacuum operations left in this table due to their visibility in transactions
<code>index_vacuum_count</code>	<code>int8</code>	Number of times indexes on this table were vacuumed
<code>rev_all_frozen_pages</code>	<code>int8</code>	

Column Type	Description
	Number of times the all-frozen mark in the visibility map was removed for pages of this table
rev_all_visible_pages int8	Number of times the all-visible mark in the visibility map was removed for pages of this table
wal_records int8	Total number of WAL records generated by vacuum operations performed on this table
wal_fpi int8	Total number of WAL full page images generated by vacuum operations performed on this table
wal_bytes numeric	Total amount of WAL bytes generated by vacuum operations performed on this table
blk_read_time int8	Time spent reading database blocks by vacuum operations performed on this table, in milliseconds (if track_io_timing is enabled, otherwise zero)
blk_write_time int8	Time spent writing database blocks by vacuum operations performed on this table, in milliseconds (if track_io_timing is enabled, otherwise zero)
delay_time float8	Time spent sleeping in a vacuum delay point by vacuum operations performed on this table, in milliseconds (see Section 19.4.4 for details)
system_time float8	System CPU time of vacuuming this table, in milliseconds
user_time float8	User CPU time of vacuuming this table, in milliseconds
total_time float8	Total time of vacuuming this table, in milliseconds
interrupts integer	Number of times vacuum operations performed on this table were interrupted on any errors

Columns `total_*`, `wal_*` and `blk_*` include data on vacuuming indexes on this table, while columns `system_time` and `user_time` only include data on vacuuming the heap.

55.34. pg_tables

The view `pg_tables` provides access to useful information about each table in the database.

Table 55.34. pg_tables Columns

Column Type	Description
schemaname name (references pg_namespace .nspname)	Name of schema containing table
tablename name (references pg_class .relname)	Name of table
tableowner name (references pg_authid .rolname)	Name of table's owner
tablespace name (references pg_tablespace .spcname)	Name of tablespace containing table (null if default for database)
hasindexes bool (references pg_class .relhasindex)	True if table has (or recently had) any indexes
hasrules bool (references pg_class .relhasrules)	

Column Type	Description
	True if table has (or once had) rules
hastriggers bool (references pg_class .relhastriggers)	True if table has (or once had) triggers
rowsecurity bool (references pg_class .relrowsecurity)	True if row security is enabled on the table

55.35. pg_timezone_abbrevs

The view `pg_timezone_abbrevs` provides a list of time zone abbreviations that are currently recognized by the datetime input routines. The contents of this view change when the [timezone_abbreviations](#) run-time parameter is modified.

Table 55.35. pg_timezone_abbrevs Columns

Column Type	Description
abbrev text	Time zone abbreviation
utc_offset interval	Offset from UTC (positive means east of Greenwich)
is_dst bool	True if this is a daylight-savings abbreviation

While most timezone abbreviations represent fixed offsets from UTC, there are some that have historically varied in value (see [Section B.4](#) for more information). In such cases this view presents their current meaning.

55.36. pg_timezone_names

The view `pg_timezone_names` provides a list of time zone names that are recognized by `SET TIMEZONE`, along with their associated abbreviations, UTC offsets, and daylight-savings status. (Technically, Postgres Pro does not use UTC because leap seconds are not handled.) Unlike the abbreviations shown in [pg_timezone_abbrevs](#), many of these names imply a set of daylight-savings transition date rules. Therefore, the associated information changes across local DST boundaries. The displayed information is computed based on the current value of `CURRENT_TIMESTAMP`.

Table 55.36. pg_timezone_names Columns

Column Type	Description
name text	Time zone name
abbrev text	Time zone abbreviation
utc_offset interval	Offset from UTC (positive means east of Greenwich)
is_dst bool	True if currently observing daylight savings

55.37. pg_user

The view `pg_user` provides access to information about database users. This is simply a publicly readable view of [pg_shadow](#) that blanks out the password field.

Table 55.37. pg_user Columns

Column Type	Description
username name	User name
usesysid oid	ID of this user
usecreatedb bool	User can create databases
usesuper bool	User is a superuser
userepl bool	User can initiate streaming replication and put the system in and out of backup mode.
usebypassrls bool	User bypasses every row-level security policy, see Section 5.9 for more information.
passwd text	Not the password (always reads as <code>*****</code>)
valuntil timestampz	Password expiry time (only used for password authentication)
useconfig text[]	Session defaults for run-time configuration variables

55.38. pg_user_mappings

The view `pg_user_mappings` provides access to information about user mappings. This is essentially a publicly readable view of `pg_user_mapping` that leaves out the `options` field if the user has no rights to use it.

Table 55.38. pg_user_mappings Columns

Column Type	Description
umid oid (references <code>pg_user_mapping</code> .oid)	OID of the user mapping
srvid oid (references <code>pg_foreign_server</code> .oid)	The OID of the foreign server that contains this mapping
srvname name (references <code>pg_foreign_server</code> .srvname)	Name of the foreign server
umuser oid (references <code>pg_authid</code> .oid)	OID of the local role being mapped, or zero if the user mapping is public
username name	Name of the local user to be mapped
umoptions text[]	User mapping specific options, as “keyword=value” strings

To protect password information stored as a user mapping option, the `umoptions` column will read as null unless one of the following applies:

- current user is the user being mapped, and owns the server or holds `USAGE` privilege on it
- current user is the server owner and mapping is for `PUBLIC`
- current user is a superuser

55.39. pg_views

The view `pg_views` provides access to useful information about each view in the database.

Table 55.39. pg_views Columns

Column Type	Description
<code>schemaname name</code> (references <code>pg_namespace .nspname</code>)	Name of schema containing view
<code>viewname name</code> (references <code>pg_class .relname</code>)	Name of view
<code>viewowner name</code> (references <code>pg_authid .rolname</code>)	Name of view's owner
<code>definition text</code>	View definition (a reconstructed <code>SELECT</code> query)

55.40. pg_wait_events

The view `pg_wait_events` provides description about the wait events.

Table 55.40. pg_wait_events Columns

Column Type	Description
<code>type text</code>	Wait event type
<code>name text</code>	Wait event name
<code>description text</code>	Wait event description

55.41. Global Views for Distributed System

Postgres Pro Shardman has a list of global views based on the corresponding local views. The definition of global view columns is the same as in its corresponding local view. Fetching from a global view returns a union of rows from the corresponding local views. The rows are fetched from each of their cluster nodes. Another difference is that the global views have an added column `rgid`. The `rgid` value shows the replication group ID of the cluster node from which a row is fetched.

55.41.1. Global Views for Statistics

Below is the list of the statistics-related global views with links to their corresponding local views:

Table 55.41. Statistics-related global and local views

Global view	Local view	Description
<code>shardman.gv_stats</code>	pg_stats	One row per planner statistics.
<code>shardman.gv_stats_ext</code>	pg_stats_ext	Provides access to information about each extended statistics object in the database.
<code>shardman.gv_stats_ext_exprs</code>	pg_stats_ext_exprs	Provides access to information about all expressions included in extended statistics objects.
<code>shardman.gv_stat_activity</code>	pg_stat_activity	One row per server process, showing information related

Global view	Local view	Description
		to the current activity of that process.
shardman.gv_stats_vacuum_database	pg_stats_vacuum_database	One row for each index in the current database (including TOAST table indexes), showing statistics about vacuuming that specific index.
shardman.gv_stats_vacuum_indexes	pg_stats_vacuum_indexes	One row for each database in the current cluster, showing statistics about vacuuming that database.
shardman.gv_stats_vacuum_tables	pg_stats_vacuum_tables	One row for each table in the current database (including TOAST tables), showing statistics about vacuuming that specific table.
shardman.gv_stat_replication	pg_stat_replication	One row per WAL sender process, showing statistics about replication to that sender's connected standby server.
shardman.gv_stat_replication_slots	pg_stat_replication_slots	One row per replication slot, showing statistics about the replication slot's usage.
shardman.gv_pgpro_stat_wal_activity	pgpro_stat_wal_activity	One row per server process, showing information related to the current activity of that process.
shardman.gv_stat_subscription	pg_stat_subscription	One row per subscription for main worker (with null PID if the worker is not running), and additional rows for workers handling the initial data copy of the subscribed tables.
shardman.gv_stat_subscription_stats	pg_stat_subscription_stats	One row per subscription.
shardman.gv_stat_ssl	pg_stat_ssl	One row per backend or WAL sender process, showing statistics about SSL usage on this connection.
shardman.gv_stat_gssapi	pg_stat_gssapi	One row per backend, showing information about GSSAPI usage on this connection.
shardman.gv_stat_archiver	pg_stat_archiver	One row only, showing statistics about the WAL archiver process's activity.
shardman.gv_stat_io	pg_stat_io	One row only for each combination of backend type, target I/O object, and I/O context, showing cluster-wide I/O statistics.

Global view	Local view	Description
shardman.gv_stat_bgwriter	pg_stat_bgwriter	One row only, showing statistics about the background writer process's activity.
shardman.gv_stat_progress_analyze	pg_stat_progress_analyze	One row for each backend (including autovacuum worker processes) running ANALYZE, showing current progress.
shardman.gv_stat_progress_basebackup	pg_stat_progress_basebackup	One row for each WAL sender process streaming a base backup, showing current progress.
shardman.gv_stat_progress_cluster	pg_stat_progress_cluster	One row for each backend running CLUSTER or VACUUM FULL, showing current progress.
shardman.gv_stat_checkpoint	pg_stat_checkpoint	One row only, containing data about the checkpoint process of the cluster.
shardman.gv_statistic_ext	pg_statistic_ext	Extended planner statistics (definition).
shardman.gv_stat_progress_create_index	pg_stat_progress_create_index	One row for each backend running CREATE INDEX or REINDEX, showing current progress.
shardman.gv_stat_progress_vacuum	pg_stat_progress_vacuum	One row for each backend (including autovacuum worker processes) that is currently vacuuming.
shardman.gv_stat_progress_copy	pg_stat_progress_copy	One row for each backend running COPY, showing current progress.
shardman.gv_stat_wal	pg_stat_wal	One row only, showing statistics about WAL activity.
shardman.gv_stat_database	pg_stat_database	One row per database, showing database-wide statistics about query cancels due to conflict with recovery on standby servers.
shardman.gv_stat_database_conflicts	pg_stat_database_conflicts	One row per database, showing database-wide statistics about query cancels occurring due to conflicts with recovery on standby servers. This view will only contain information on standby servers, since conflicts do not occur on primary servers.
shardman.gv_stat_all_tables	pg_stat_all_tables	One row for each table in the current database, showing statistics about accesses to that specific table.
shardman.gv_stat_sys_tables	pg_stat_sys_tables	Same as pg_stat_sys_tables , only as a global view.

Global view	Local view	Description
shardman.gv_stat_user_tables	pg_stat_user_tables	Same as pg_stat_user_tables , only as a global view.
shardman.gv_stat_all_indexes	pg_stat_all_indexes	One row for each index in the current database, showing statistics about accesses to that specific index.
shardman.gv_stat_sys_indexes	pg_stat_sys_indexes	Same as pg_stat_sys_indexes , only as a global view.
shardman.gv_stat_user_indexes	pg_stat_user_indexes	Same as pg_stat_user_indexes , only as a global view.
shardman.gv_statio_all_tables	pg_statio_all_tables	One row for each table in the current database, showing statistics about I/O on that specific table.
shardman.gv_statio_sys_tables	pg_statio_sys_tables	Same as pg_statio_sys_tables , only as a global view.
shardman.gv_statio_user_tables	pg_statio_user_tables	Same as pg_statio_user_tables , only as a global view.
shardman.gv_statio_all_indexes	pg_statio_all_indexes	One row for each index in the current database, showing statistics about I/O on that specific index.
shardman.gv_statio_sys_indexes	pg_statio_sys_indexes	Same as pg_statio_sys_indexes , only as a global view.
shardman.gv_statio_user_indexes	pg_statio_user_indexes	Same as pg_statio_user_indexes , only as a global view.
shardman.gv_statio_all_sequences	pg_statio_all_sequences	One row for each sequence in the current database, showing statistics about I/O on that specific sequence.
shardman.gv_statio_sys_sequences	pg_statio_sys_sequences	Same as pg_statio_sys_sequences , only as a global view.
shardman.gv_statio_user_sequences	pg_statio_user_sequences	Same as pg_statio_user_sequences , only as a global view.
shardman.gv_stat_user_functions	pg_stat_user_functions	One row for each tracked function, showing statistics about executions of that function.
shardman.gv_stat_slru	pg_stat_slru	One row per SLRU, showing statistics of operations.
shardman.gv_stat_csn	shardman.pg_stat_csn	One row showing statistics on delays that take place during import of CSN snapshots.
shardman.gv_stat_monitor	shardman.pg_stat_monitor	One row showing metrics of the Postgres Pro Shardman monitor.
shardman.gv_stat_netusage	shardman.pg_stat_net_usage	One row showing the cumulative network traffic between Postgres Pro Shardman cluster nodes.

Global view	Local view	Description
<code>shardman.gv_stat_xact_time</code>	<code>shardman.pg_stat_xact_time</code>	One row showing statistics for the time spent on a transaction.
<code>shardman.gv_lock_graph</code>	<code>shardman.lock_graph</code>	One row showing a graph of locks between processes on Postgres Pro Shardman cluster nodes including external locks. This view is based on the <code>pg_locks</code> and <code>pg_prepared_xacts</code> system views and on the <code>pg_stat_activity</code> view of the cumulative statistics system .
<code>shardman.gv_stat_foreign_bytes</code>	<code>shardman.pg_stat_foreign_stat_bytes</code>	One row showing the amount of statistics for foreign relations transferred over the network between Postgres Pro Shardman cluster nodes.
<code>shardman.gv_stat_recovery_prefetch</code>	<code>pg_stat_recovery_prefetch</code>	One row, the columns <code>wal_distance</code> , <code>block_distance</code> and <code>io_depth</code> show current values, and the other columns show cumulative counters that can be reset with the <code>pg_stat_reset_shared</code> function.
<code>shardman.gv_stat_wal_receiver</code>	<code>pg_stat_wal_receiver</code>	One row, showing statistics about the WAL receiver from that receiver's connected server.
<code>shardman.gv_stat_xact_all_tables</code>	<code>pg_stat_xact_all_tables</code>	Similar to <code>pg_stat_all_tables</code> , but counts actions taken so far within the current transaction (which are <i>not</i> yet included in <code>pg_stat_all_tables</code> and related views). The columns for numbers of live and dead rows and vacuum and analyze actions are not present in this view.
<code>shardman.gv_stat_xact_sys_tables</code>	<code>pg_stat_xact_sys_tables</code>	Same as <code>pg_stat_xact_sys_tables</code> , only as a global view.
<code>shardman.gv_stat_xact_user_tables</code>	<code>pg_stat_xact_user_tables</code>	Same as <code>pg_stat_xact_user_tables</code> , only as a global view
<code>shardman.gv_stat_xact_user_functions</code>	<code>pg_stat_xact_user_functions</code>	Similar to <code>pg_stat_user_functions</code> , but counts only calls during the current transaction (which are <i>not</i> yet included in <code>pg_stat_user_functions</code>).
<code>shardman.gv_sharded_tables</code>		Displays information on all the sharded tables in the cluster.
<code>shardman.gv_global_tables</code>		Displays information on all the global tables in the cluster.

55.41.2. Global Views for System Catalog and System Views

Below is the list of the global views that relate to the system catalog and system views, and links to their corresponding local views:

Table 55.42. Global and local views for system catalog and system views

Global view	Local view	Description
shardman.gv_aggregate	pg_aggregate	Stores information about aggregate functions
shardman.gv_am	pg_am	Relation access methods
shardman.gv_amop	pg_amop	Access method operators
shardman.gv_amproc	pg_amproc	Access method support functions
shardman.gv_attrdef	pg_attrdef	Column default values
shardman.gv_attribute	pg_attribute	Table columns (“attributes”)
shardman.gv_auth_members	pg_auth_members	Authorization identifier membership relationships
shardman.gv_available_extension_versions	pg_available_extension_versions	Specific extension versions that are available for installation
shardman.gv_available_extensions	pg_available_extensions	Extensions that are available for installation
shardman.gv_cast	pg_cast	Casts (data type conversions)
shardman.gv_class	pg_class	Tables, indexes, sequences, views (“relations”)
shardman.gv_collation	pg_collation	Collations (locale information)
shardman.gv_config	pg_config	Compile-time configuration parameters of the currently installed version of Postgres Pro
shardman.gv_constraint	pg_constraint	Check constraints, unique constraints, primary key constraints, foreign key constraints
shardman.gv_conversion	pg_conversion	Encoding conversion information
shardman.gv_database	pg_database	Databases within this database cluster
shardman.gv_db_role_setting	pg_db_role_setting	Per-role and per-database settings
shardman.gv_default_acl	pg_default_acl	Default privileges for object types
shardman.gv_depend	pg_depend	Dependencies between database objects
shardman.gv_description	pg_description	Descriptions or comments on database objects
shardman.gv_enum	pg_enum	Enum label and value definitions
shardman.gv_event_trigger	pg_event_trigger	Event triggers
shardman.gv_extension	pg_extension	Installed extensions
shardman.gv_file_settings	pg_file_settings	Installed extensions
shardman.gv_foreign_data_wrapper	pg_foreign_data_wrapper	Foreign-data wrapper definitions

Global view	Local view	Description
shardman.gv_foreign_server	pg_foreign_server	Foreign server definitions
shardman.gv_foreign_table	pg_foreign_table	Additional foreign table information
shardman.gv_group	pg_group	Exists for backwards compatibility: it emulates a catalog that existed in PostgreSQL before version 8.1
shardman.gv_hba_file_rules	pg_hba_file_rules	Summary of the contents of the client authentication configuration file <code>pg_hba.conf</code>
shardman.gv_ident_file_mappings	pg_ident_file_mappings	Summary of the contents of the client user name mapping configuration file <code>pg_hba.conf</code>
shardman.gv_index	pg_index	Additional index information
shardman.gv_indexes	pg_indexes	Provides access to useful information about each index in the database
shardman.gv_inherits	pg_inherits	Table inheritance hierarchy
shardman.gv_init_privs	pg_init_privs	Object initial privileges
shardman.gv_language	pg_language	Languages for writing functions
shardman.gv_largeobject	pg_largeobject	Data pages for large objects
shardman.gv_largeobject_metadata	pg_largeobject_metadata	Metadata associated with large objects
shardman.gv_matviews	pg_matviews	Provides access to useful information about each materialized view in the database
shardman.gv_namespace	pg_namespace	Schemas
shardman.gv_opclass	pg_opclass	Access method operator classes
shardman.gv_operator	pg_operator	Operators
shardman.gv_opfamily	pg_opfamily	Access method operator families
shardman.gv_parameter_acl	pg_parameter_acl	Configuration parameters for which privileges have been granted
shardman.gv_partitioned_table	pg_partitioned_table	Information about partition key of tables
shardman.gv_proc	pg_proc	Functions and procedures
shardman.gv_profile	pg_profile	Profiles, a set of authentication restrictions
shardman.gv_publication	pg_publication	Publications for logical replication
shardman.gv_publication_namespace	pg_publication_namespace	Schema to publication mapping
shardman.gv_publication_rel	pg_publication_rel	Relation to publication mapping

Global view	Local view	Description
shardman.gv_publication_tables	pg_publication_tables	Information about the mapping between publications and information of tables they contain
shardman.gv_range	pg_range	Information about range types
shardman.gv_replication_origin	pg_replication_origin	Registered replication origins
shardman.gv_replication_origin_status	pg_replication_origin_status	Information about how far replay for a certain origin has progressed
shardman.gv_replication_slots	pg_replication_slots	Provides a listing of all replication slots that currently exist on the database cluster, along with their current state
shardman.gv_rewrite	pg_rewrite	Query rewrite rules
shardman.gv_rules	pg_rules	Provides access to useful information about query rewrite rules
shardman.gv_seclabel	pg_seclabel	Security labels on database objects
shardman.gv_seclabels	pg_seclabels	Provides information about security labels
shardman.gv_sequence	pg_sequence	Information about sequences
shardman.gv_sequences	pg_sequences	Provides access to useful information about each sequence in the database
shardman.gv_settings	pg_settings	Provides access to run-time parameters of the server
shardman.gv_shdepend	pg_shdepend	Dependencies on shared objects
shardman.gv_shdescription	pg_shdescription	Comments on shared objects
shardman.gv_shseclabel	pg_shseclabel	Security labels on shared database objects
shardman.gv_subscription	pg_subscription	Logical replication subscriptions
shardman.gv_subscription_rel	pg_subscription_rel	Relation state for subscriptions
shardman.gv_tablespace	pg_tablespace	Tablespaces within this database cluster
shardman.gv_tables	pg_tables	Provides access to useful information about each table in the database
shardman.gv_prepared_xacts	pg_prepared_xacts	Provides information about transactions that are currently prepared for two-phase commit
shardman.gv_timezone_names	pg_timezone_names	List of time zone names that are recognized by <code>SET TIMEZONE</code> , along with their associated abbreviations, UTC offsets, and daylight-savings status

Global view	Local view	Description
shardman.gv_timezone_abbrevs	pg_timezone_abbrevs	List of time zone abbreviations that are currently recognized by the datetime input routines
shardman.gv_transform	pg_transform	Transforms (data type to procedural language conversions)
shardman.gv_trigger	pg_trigger	Triggers
shardman.gv_ts_config	pg_ts_config	Text search configurations
shardman.gv_ts_config_map	pg_ts_config_map	Text search configurations' token mappings
shardman.gv_ts_dict	pg_ts_dict	Text search dictionaries
shardman.gv_ts_parser	pg_ts_parser	Text search parsers
shardman.gv_ts_template	pg_ts_template	Text search templates
shardman.gv_type	pg_type	Data types
shardman.gv_user_mapping	pg_user_mapping	Mappings of users to foreign servers
shardman.gv_user_mappings	pg_user_mappings	Provides access to information about user mappings
shardman.gv_views	pg_views	Provides access to useful information about each view in the database
shardman.gv_wait_events	pg_wait_events	Provides description about the wait events
shardman.gv_locks	pg_locks	Provides access to information about the locks held by active processes within the database server.
shardman.gv_shmem_allocations	pg_shmem_allocations	Shows allocations made from the server's main shared memory segment.

Chapter 56. Frontend/Backend Protocol

Postgres Pro uses a message-based protocol for communication between frontends and backends (clients and servers). The protocol is supported over TCP/IP and also over Unix-domain sockets. Port number 5432 has been registered with IANA as the customary TCP port number for servers supporting this protocol, but in practice any non-privileged port number can be used.

This document describes version 3.0 of the protocol, implemented in PostgreSQL 7.4 and later. For descriptions of the earlier protocol versions, see previous releases of the PostgreSQL documentation. A single server can support multiple protocol versions. The initial startup-request message tells the server which protocol version the client is attempting to use. If the major version requested by the client is not supported by the server, the connection will be rejected (for example, this would occur if the client requested protocol version 4.0, which does not exist as of this writing). If the minor version requested by the client is not supported by the server (e.g., the client requests version 3.1, but the server supports only 3.0), the server may either reject the connection or may respond with a `NegotiateProtocolVersion` message containing the highest minor protocol version which it supports. The client may then choose either to continue with the connection using the specified protocol version or to abort the connection.

In order to serve multiple clients efficiently, the server launches a new “backend” process for each client. In the current implementation, a new child process is created immediately after an incoming connection is detected. This is transparent to the protocol, however. For purposes of the protocol, the terms “backend” and “server” are interchangeable; likewise “frontend” and “client” are interchangeable.

56.1. Overview

The protocol has separate phases for startup and normal operation. In the startup phase, the frontend opens a connection to the server and authenticates itself to the satisfaction of the server. (This might involve a single message, or multiple messages depending on the authentication method being used.) If all goes well, the server then sends status information to the frontend, and finally enters normal operation. Except for the initial startup-request message, this part of the protocol is driven by the server.

During normal operation, the frontend sends queries and other commands to the backend, and the backend sends back query results and other responses. There are a few cases (such as `NOTIFY`) wherein the backend will send unsolicited messages, but for the most part this portion of a session is driven by frontend requests.

Termination of the session is normally by frontend choice, but can be forced by the backend in certain cases. In any case, when the backend closes the connection, it will roll back any open (incomplete) transaction before exiting.

Within normal operation, SQL commands can be executed through either of two sub-protocols. In the “simple query” protocol, the frontend just sends a textual query string, which is parsed and immediately executed by the backend. In the “extended query” protocol, processing of queries is separated into multiple steps: parsing, binding of parameter values, and execution. This offers flexibility and performance benefits, at the cost of extra complexity.

Normal operation has additional sub-protocols for special operations such as `COPY`.

56.1.1. Messaging Overview

All communication is through a stream of messages. The first byte of a message identifies the message type, and the next four bytes specify the length of the rest of the message (this length count includes itself, but not the message-type byte). The remaining contents of the message are determined by the message type. For historical reasons, the very first message sent by the client (the startup message) has no initial message-type byte.

To avoid losing synchronization with the message stream, both servers and clients typically read an entire message into a buffer (using the byte count) before attempting to process its contents. This allows easy recovery if an error is detected while processing the contents. In extreme situations (such as not

having enough memory to buffer the message), the receiver can use the byte count to determine how much input to skip before it resumes reading messages.

Conversely, both servers and clients must take care never to send an incomplete message. This is commonly done by marshaling the entire message in a buffer before beginning to send it. If a communications failure occurs partway through sending or receiving a message, the only sensible response is to abandon the connection, since there is little hope of recovering message-boundary synchronization.

56.1.2. Extended Query Overview

In the extended-query protocol, execution of SQL commands is divided into multiple steps. The state retained between steps is represented by two types of objects: *prepared statements* and *portals*. A prepared statement represents the result of parsing and semantic analysis of a textual query string. A prepared statement is not in itself ready to execute, because it might lack specific values for *parameters*. A portal represents a ready-to-execute or already-partially-executed statement, with any missing parameter values filled in. (For `SELECT` statements, a portal is equivalent to an open cursor, but we choose to use a different term since cursors don't handle non-`SELECT` statements.)

The overall execution cycle consists of a *parse* step, which creates a prepared statement from a textual query string; a *bind* step, which creates a portal given a prepared statement and values for any needed parameters; and an *execute* step that runs a portal's query. In the case of a query that returns rows (`SELECT`, `SHOW`, etc.), the execute step can be told to fetch only a limited number of rows, so that multiple execute steps might be needed to complete the operation.

The backend can keep track of multiple prepared statements and portals (but note that these exist only within a session, and are never shared across sessions). Existing prepared statements and portals are referenced by names assigned when they were created. In addition, an “unnamed” prepared statement and portal exist. Although these behave largely the same as named objects, operations on them are optimized for the case of executing a query only once and then discarding it, whereas operations on named objects are optimized on the expectation of multiple uses.

56.1.3. Formats and Format Codes

Data of a particular data type might be transmitted in any of several different *formats*. As of PostgreSQL 7.4 the only supported formats are “text” and “binary”, but the protocol makes provision for future extensions. The desired format for any value is specified by a *format code*. Clients can specify a format code for each transmitted parameter value and for each column of a query result. Text has format code zero, binary has format code one, and all other format codes are reserved for future definition.

The text representation of values is whatever strings are produced and accepted by the input/output conversion functions for the particular data type. In the transmitted representation, there is no trailing null character; the frontend must add one to received values if it wants to process them as C strings. (The text format does not allow embedded nulls, by the way.)

Binary representations for integers use network byte order (most significant byte first). For other data types consult the documentation to learn about the binary representation. Keep in mind that binary representations for complex data types might change across server versions; the text format is usually the more portable choice.

56.2. Message Flow

This section describes the message flow and the semantics of each message type. (Details of the exact representation of each message appear in [Section 56.7.](#)) There are several different sub-protocols depending on the state of the connection: start-up, query, function call, `COPY`, and termination. There are also special provisions for asynchronous operations (including notification responses and command cancellation), which can occur at any time after the start-up phase.

56.2.1. Start-up

To begin a session, a frontend opens a connection to the server and sends a startup message. This message includes the names of the user and of the database the user wants to connect to; it also identifies the

particular protocol version to be used. (Optionally, the startup message can include additional settings for run-time parameters.) The server then uses this information and the contents of its configuration files (such as `pg_hba.conf`) to determine whether the connection is provisionally acceptable, and what additional authentication is required (if any).

The server then sends an appropriate authentication request message, to which the frontend must reply with an appropriate authentication response message (such as a password). For all authentication methods except GSSAPI, SSPI and SASL, there is at most one request and one response. In some methods, no response at all is needed from the frontend, and so no authentication request occurs. For GSSAPI, SSPI and SASL, multiple exchanges of packets may be needed to complete the authentication.

The authentication cycle ends with the server either rejecting the connection attempt (`ErrorResponse`), or sending `AuthenticationOk`.

The possible messages from the server in this phase are:

`ErrorResponse`

The connection attempt has been rejected. The server then immediately closes the connection.

`AuthenticationOk`

The authentication exchange is successfully completed.

`AuthenticationKerberosV5`

The frontend must now take part in a Kerberos V5 authentication dialog (not described here, part of the Kerberos specification) with the server. If this is successful, the server responds with an `AuthenticationOk`, otherwise it responds with an `ErrorResponse`. This is no longer supported.

`AuthenticationCleartextPassword`

The frontend must now send a `PasswordMessage` containing the password in clear-text form. If this is the correct password, the server responds with an `AuthenticationOk`, otherwise it responds with an `ErrorResponse`.

`AuthenticationMD5Password`

The frontend must now send a `PasswordMessage` containing the password (with user name) encrypted via MD5, then encrypted again using the 4-byte random salt specified in the `AuthenticationMD5Password` message. If this is the correct password, the server responds with an `AuthenticationOk`, otherwise it responds with an `ErrorResponse`. The actual `PasswordMessage` can be computed in SQL as `concat('md5', md5(concat(md5(concat(password, username)), random-salt)))`. (Keep in mind the `md5()` function returns its result as a hex string.)

`AuthenticationGSS`

The frontend must now initiate a GSSAPI negotiation. The frontend will send a `GSSResponse` message with the first part of the GSSAPI data stream in response to this. If further messages are needed, the server will respond with `AuthenticationGSSContinue`.

`AuthenticationSSPI`

The frontend must now initiate an SSPI negotiation. The frontend will send a `GSSResponse` with the first part of the SSPI data stream in response to this. If further messages are needed, the server will respond with `AuthenticationGSSContinue`.

`AuthenticationGSSContinue`

This message contains the response data from the previous step of GSSAPI or SSPI negotiation (`AuthenticationGSS`, `AuthenticationSSPI` or a previous `AuthenticationGSSContinue`). If the GSSAPI or SSPI data in this message indicates more data is needed to complete the authentication, the frontend must send that data as another `GSSResponse` message. If GSSAPI or SSPI authentication is completed by this message, the server will next send `AuthenticationOk` to indicate successful authentication or `ErrorResponse` to indicate failure.

AuthenticationSASL

The frontend must now initiate a SASL negotiation, using one of the SASL mechanisms listed in the message. The frontend will send a `SASLInitialResponse` with the name of the selected mechanism, and the first part of the SASL data stream in response to this. If further messages are needed, the server will respond with `AuthenticationSASLContinue`. See [Section 56.3](#) for details.

AuthenticationSASLContinue

This message contains challenge data from the previous step of SASL negotiation (`AuthenticationSASL`, or a previous `AuthenticationSASLContinue`). The frontend must respond with a `SASLResponse` message.

AuthenticationSASLFinal

SASL authentication has completed with additional mechanism-specific data for the client. The server will next send `AuthenticationOk` to indicate successful authentication, or an `ErrorResponse` to indicate failure. This message is sent only if the SASL mechanism specifies additional data to be sent from server to client at completion.

NegotiateProtocolVersion

The server does not support the minor protocol version requested by the client, but does support an earlier version of the protocol; this message indicates the highest supported minor version. This message will also be sent if the client requested unsupported protocol options (i.e., beginning with `_pq_`) in the startup packet. This message will be followed by an `ErrorResponse` or a message indicating the success or failure of authentication.

If the frontend does not support the authentication method requested by the server, then it should immediately close the connection.

After having received `AuthenticationOk`, the frontend must wait for further messages from the server. In this phase a backend process is being started, and the frontend is just an interested bystander. It is still possible for the startup attempt to fail (`ErrorResponse`) or the server to decline support for the requested minor protocol version (`NegotiateProtocolVersion`), but in the normal case the backend will send some `ParameterStatus` messages, `BackendKeyData`, and finally `ReadyForQuery`.

During this phase the backend will attempt to apply any additional run-time parameter settings that were given in the startup message. If successful, these values become session defaults. An error causes `ErrorResponse` and exit.

The possible messages from the backend in this phase are:

BackendKeyData

This message provides secret-key data that the frontend must save if it wants to be able to issue cancel requests later. The frontend should not respond to this message, but should continue listening for a `ReadyForQuery` message.

ParameterStatus

This message informs the frontend about the current (initial) setting of backend parameters, such as [client encoding](#) or [DateStyle](#). The frontend can ignore this message, or record the settings for its future use; see [Section 56.2.7](#) for more details. The frontend should not respond to this message, but should continue listening for a `ReadyForQuery` message.

ReadyForQuery

Start-up is completed. The frontend can now issue commands.

ErrorResponse

Start-up failed. The connection is closed after sending this message.

NoticeResponse

A warning message has been issued. The frontend should display the message but continue listening for ReadyForQuery or ErrorResponse.

The ReadyForQuery message is the same one that the backend will issue after each command cycle. Depending on the coding needs of the frontend, it is reasonable to consider ReadyForQuery as starting a command cycle, or to consider ReadyForQuery as ending the start-up phase and each subsequent command cycle.

56.2.2. Simple Query

A simple query cycle is initiated by the frontend sending a Query message to the backend. The message includes an SQL command (or commands) expressed as a text string. The backend then sends one or more response messages depending on the contents of the query command string, and finally a ReadyForQuery response message. ReadyForQuery informs the frontend that it can safely send a new command. (It is not actually necessary for the frontend to wait for ReadyForQuery before issuing another command, but the frontend must then take responsibility for figuring out what happens if the earlier command fails and already-issued later commands succeed.)

The possible response messages from the backend are:

CommandComplete

An SQL command completed normally.

CopyInResponse

The backend is ready to copy data from the frontend to a table; see [Section 56.2.6](#).

CopyOutResponse

The backend is ready to copy data from a table to the frontend; see [Section 56.2.6](#).

RowDescription

Indicates that rows are about to be returned in response to a `SELECT`, `FETCH`, etc. query. The contents of this message describe the column layout of the rows. This will be followed by a DataRow message for each row being returned to the frontend.

DataRow

One of the set of rows returned by a `SELECT`, `FETCH`, etc. query.

EmptyQueryResponse

An empty query string was recognized.

ErrorResponse

An error has occurred.

ReadyForQuery

Processing of the query string is complete. A separate message is sent to indicate this because the query string might contain multiple SQL commands. (CommandComplete marks the end of processing one SQL command, not the whole string.) ReadyForQuery will always be sent, whether processing terminates successfully or with an error.

NoticeResponse

A warning message has been issued in relation to the query. Notices are in addition to other responses, i.e., the backend will continue processing the command.

The response to a `SELECT` query (or other queries that return row sets, such as `EXPLAIN` or `SHOW`) normally consists of RowDescription, zero or more DataRow messages, and then CommandComplete. `COPY` to or

from the frontend invokes special protocol as described in [Section 56.2.6](#). All other query types normally produce only a `CommandComplete` message.

Since a query string could contain several queries (separated by semicolons), there might be several such response sequences before the backend finishes processing the query string. `ReadyForQuery` is issued when the entire string has been processed and the backend is ready to accept a new query string.

If a completely empty (no contents other than whitespace) query string is received, the response is `EmptyQueryResponse` followed by `ReadyForQuery`.

In the event of an error, `ErrorResponse` is issued followed by `ReadyForQuery`. All further processing of the query string is aborted by `ErrorResponse` (even if more queries remained in it). Note that this might occur partway through the sequence of messages generated by an individual query.

In simple Query mode, the format of retrieved values is always text, except when the given command is a `FETCH` from a cursor declared with the `BINARY` option. In that case, the retrieved values are in binary format. The format codes given in the `RowDescription` message tell which format is being used.

A frontend must be prepared to accept `ErrorResponse` and `NoticeResponse` messages whenever it is expecting any other type of message. See also [Section 56.2.7](#) concerning messages that the backend might generate due to outside events.

Recommended practice is to code frontends in a state-machine style that will accept any message type at any time that it could make sense, rather than wiring in assumptions about the exact sequence of messages.

56.2.2.1. Multiple Statements in a Simple Query

When a simple Query message contains more than one SQL statement (separated by semicolons), those statements are executed as a single transaction, unless explicit transaction control commands are included to force a different behavior. For example, if the message contains

```
INSERT INTO mytable VALUES(1);
SELECT 1/0;
INSERT INTO mytable VALUES(2);
```

then the divide-by-zero failure in the `SELECT` will force rollback of the first `INSERT`. Furthermore, because execution of the message is abandoned at the first error, the second `INSERT` is never attempted at all.

If instead the message contains

```
BEGIN;
INSERT INTO mytable VALUES(1);
COMMIT;
INSERT INTO mytable VALUES(2);
SELECT 1/0;
```

then the first `INSERT` is committed by the explicit `COMMIT` command. The second `INSERT` and the `SELECT` are still treated as a single transaction, so that the divide-by-zero failure will roll back the second `INSERT`, but not the first one.

This behavior is implemented by running the statements in a multi-statement Query message in an *implicit transaction block* unless there is some explicit transaction block for them to run in. The main difference between an implicit transaction block and a regular one is that an implicit block is closed automatically at the end of the Query message, either by an implicit commit if there was no error, or an implicit rollback if there was an error. This is similar to the implicit commit or rollback that happens for a statement executed by itself (when not in a transaction block).

If the session is already in a transaction block, as a result of a `BEGIN` in some previous message, then the Query message simply continues that transaction block, whether the message contains one statement or several. However, if the Query message contains a `COMMIT` or `ROLLBACK` closing the existing transaction

block, then any following statements are executed in an implicit transaction block. Conversely, if a `BEGIN` appears in a multi-statement Query message, then it starts a regular transaction block that will only be terminated by an explicit `COMMIT` or `ROLLBACK`, whether that appears in this Query message or a later one. If the `BEGIN` follows some statements that were executed as an implicit transaction block, those statements are not immediately committed; in effect, they are retroactively included into the new regular transaction block.

A `COMMIT` or `ROLLBACK` appearing in an implicit transaction block is executed as normal, closing the implicit block; however, a warning will be issued since a `COMMIT` or `ROLLBACK` without a previous `BEGIN` might represent a mistake. If more statements follow, a new implicit transaction block will be started for them.

Savepoints are not allowed in an implicit transaction block, since they would conflict with the behavior of automatically closing the block upon any error.

Remember that, regardless of any transaction control commands that may be present, execution of the Query message stops at the first error. Thus for example given

```
BEGIN;  
SELECT 1/0;  
ROLLBACK;
```

in a single Query message, the session will be left inside a failed regular transaction block, since the `ROLLBACK` is not reached after the divide-by-zero error. Another `ROLLBACK` will be needed to restore the session to a usable state.

Another behavior of note is that initial lexical and syntactic analysis is done on the entire query string before any of it is executed. Thus simple errors (such as a misspelled keyword) in later statements can prevent execution of any of the statements. This is normally invisible to users since the statements would all roll back anyway when done as an implicit transaction block. However, it can be visible when attempting to do multiple transactions within a multi-statement Query. For instance, if a typo turned our previous example into

```
BEGIN;  
INSERT INTO mytable VALUES(1);  
COMMIT;  
INSERT INTO mytable VALUES(2);  
SELC 1/0;
```

then none of the statements would get run, resulting in the visible difference that the first `INSERT` is not committed. Errors detected at semantic analysis or later, such as a misspelled table or column name, do not have this effect.

56.2.3. Extended Query

The extended query protocol breaks down the above-described simple query protocol into multiple steps. The results of preparatory steps can be re-used multiple times for improved efficiency. Furthermore, additional features are available, such as the possibility of supplying data values as separate parameters instead of having to insert them directly into a query string.

In the extended protocol, the frontend first sends a Parse message, which contains a textual query string, optionally some information about data types of parameter placeholders, and the name of a destination prepared-statement object (an empty string selects the unnamed prepared statement). The response is either ParseComplete or ErrorResponse. Parameter data types can be specified by OID; if not given, the parser attempts to infer the data types in the same way as it would do for untyped literal string constants.

Note

A parameter data type can be left unspecified by setting it to zero, or by making the array of parameter type OIDs shorter than the number of parameter symbols ($\$n$) used in the query string.

Another special case is that a parameter's type can be specified as `void` (that is, the OID of the `void` pseudo-type). This is meant to allow parameter symbols to be used for function parameters that are actually OUT parameters. Ordinarily there is no context in which a `void` parameter could be used, but if such a parameter symbol appears in a function's parameter list, it is effectively ignored. For example, a function call such as `foo($1,$2,$3,$4)` could match a function with two IN and two OUT arguments, if `$3` and `$4` are specified as having type `void`.

Note

The query string contained in a Parse message cannot include more than one SQL statement; else a syntax error is reported. This restriction does not exist in the simple-query protocol, but it does exist in the extended protocol, because allowing prepared statements or portals to contain multiple commands would complicate the protocol unduly.

If successfully created, a named prepared-statement object lasts till the end of the current session, unless explicitly destroyed. An unnamed prepared statement lasts only until the next Parse statement specifying the unnamed statement as destination is issued. (Note that a simple Query message also destroys the unnamed statement.) Named prepared statements must be explicitly closed before they can be redefined by another Parse message, but this is not required for the unnamed statement. Named prepared statements can also be created and accessed at the SQL command level, using `PREPARE` and `EXECUTE`.

Once a prepared statement exists, it can be readied for execution using a Bind message. The Bind message gives the name of the source prepared statement (empty string denotes the unnamed prepared statement), the name of the destination portal (empty string denotes the unnamed portal), and the values to use for any parameter placeholders present in the prepared statement. The supplied parameter set must match those needed by the prepared statement. (If you declared any `void` parameters in the Parse message, pass NULL values for them in the Bind message.) Bind also specifies the format to use for any data returned by the query; the format can be specified overall, or per-column. The response is either `BindComplete` or `ErrorResponse`.

Note

The choice between text and binary output is determined by the format codes given in Bind, regardless of the SQL command involved. The `BINARY` attribute in cursor declarations is irrelevant when using extended query protocol.

Query planning typically occurs when the Bind message is processed. If the prepared statement has no parameters, or is executed repeatedly, the server might save the created plan and re-use it during subsequent Bind messages for the same prepared statement. However, it will do so only if it finds that a generic plan can be created that is not much less efficient than a plan that depends on the specific parameter values supplied. This happens transparently so far as the protocol is concerned.

If successfully created, a named portal object lasts till the end of the current transaction, unless explicitly destroyed. An unnamed portal is destroyed at the end of the transaction, or as soon as the next Bind statement specifying the unnamed portal as destination is issued. (Note that a simple Query message also destroys the unnamed portal.) Named portals must be explicitly closed before they can be redefined by another Bind message, but this is not required for the unnamed portal. Named portals can also be created and accessed at the SQL command level, using `DECLARE CURSOR` and `FETCH`.

Once a portal exists, it can be executed using an Execute message. The Execute message specifies the portal name (empty string denotes the unnamed portal) and a maximum result-row count (zero meaning "fetch all rows"). The result-row count is only meaningful for portals containing commands that return row sets; in other cases the command is always executed to completion, and the row count is ignored.

The possible responses to `Execute` are the same as those described above for queries issued via simple query protocol, except that `Execute` doesn't cause `ReadyForQuery` or `RowDescription` to be issued.

If `Execute` terminates before completing the execution of a portal (due to reaching a nonzero result-row count), it will send a `PortalSuspended` message; the appearance of this message tells the frontend that another `Execute` should be issued against the same portal to complete the operation. The `CommandComplete` message indicating completion of the source SQL command is not sent until the portal's execution is completed. Therefore, an `Execute` phase is always terminated by the appearance of exactly one of these messages: `CommandComplete`, `EmptyQueryResponse` (if the portal was created from an empty query string), `ErrorResponse`, or `PortalSuspended`.

At completion of each series of extended-query messages, the frontend should issue a `Sync` message. This parameterless message causes the backend to close the current transaction if it's not inside a `BEGIN/COMMIT` transaction block ("close" meaning to commit if no error, or roll back if error). Then a `ReadyForQuery` response is issued. The purpose of `Sync` is to provide a resynchronization point for error recovery. When an error is detected while processing any extended-query message, the backend issues `ErrorResponse`, then reads and discards messages until a `Sync` is reached, then issues `ReadyForQuery` and returns to normal message processing. (But note that no skipping occurs if an error is detected *while* processing `Sync` — this ensures that there is one and only one `ReadyForQuery` sent for each `Sync`.)

Note

`Sync` does not cause a transaction block opened with `BEGIN` to be closed. It is possible to detect this situation since the `ReadyForQuery` message includes transaction status information.

In addition to these fundamental, required operations, there are several optional operations that can be used with extended-query protocol.

The `Describe` message (portal variant) specifies the name of an existing portal (or an empty string for the unnamed portal). The response is a `RowDescription` message describing the rows that will be returned by executing the portal; or a `NoData` message if the portal does not contain a query that will return rows; or `ErrorResponse` if there is no such portal.

The `Describe` message (statement variant) specifies the name of an existing prepared statement (or an empty string for the unnamed prepared statement). The response is a `ParameterDescription` message describing the parameters needed by the statement, followed by a `RowDescription` message describing the rows that will be returned when the statement is eventually executed (or a `NoData` message if the statement will not return rows). `ErrorResponse` is issued if there is no such prepared statement. Note that since `Bind` has not yet been issued, the formats to be used for returned columns are not yet known to the backend; the format code fields in the `RowDescription` message will be zeroes in this case.

Tip

In most scenarios the frontend should issue one or the other variant of `Describe` before issuing `Execute`, to ensure that it knows how to interpret the results it will get back.

The `Close` message closes an existing prepared statement or portal and releases resources. It is not an error to issue `Close` against a nonexistent statement or portal name. The response is normally `CloseComplete`, but could be `ErrorResponse` if some difficulty is encountered while releasing resources. Note that closing a prepared statement implicitly closes any open portals that were constructed from that statement.

The `Flush` message does not cause any specific output to be generated, but forces the backend to deliver any data pending in its output buffers. A `Flush` must be sent after any extended-query command except `Sync`, if the frontend wishes to examine the results of that command before issuing more commands.

Without Flush, messages returned by the backend will be combined into the minimum possible number of packets to minimize network overhead.

Note

The simple Query message is approximately equivalent to the series Parse, Bind, portal Describe, Execute, Close, Sync, using the unnamed prepared statement and portal objects and no parameters. One difference is that it will accept multiple SQL statements in the query string, automatically performing the bind/describe/execute sequence for each one in succession. Another difference is that it will not return ParseComplete, BindComplete, CloseComplete, or NoData messages.

56.2.4. Pipelining

Use of the extended query protocol allows *pipelining*, which means sending a series of queries without waiting for earlier ones to complete. This reduces the number of network round trips needed to complete a given series of operations. However, the user must carefully consider the required behavior if one of the steps fails, since later queries will already be in flight to the server.

One way to deal with that is to make the whole query series be a single transaction, that is wrap it in `BEGIN ... COMMIT`. However, this does not help if one wishes for some of the commands to commit independently of others.

The extended query protocol provides another way to manage this concern, which is to omit sending Sync messages between steps that are dependent. Since, after an error, the backend will skip command messages until it finds Sync, this allows later commands in a pipeline to be skipped automatically when an earlier one fails, without the client having to manage that explicitly with `BEGIN` and `COMMIT`. Independently-committable segments of the pipeline can be separated by Sync messages.

If the client has not issued an explicit `BEGIN`, then each Sync ordinarily causes an implicit `COMMIT` if the preceding step(s) succeeded, or an implicit `ROLLBACK` if they failed. However, there are a few DDL commands (such as `CREATE DATABASE`) that cannot be executed inside a transaction block. If one of these is executed in a pipeline, it will fail unless it is the first command in the pipeline. Furthermore, upon success it will force an immediate commit to preserve database consistency. Thus a Sync immediately following one of these commands has no effect except to respond with ReadyForQuery.

When using this method, completion of the pipeline must be determined by counting ReadyForQuery messages and waiting for that to reach the number of Syncs sent. Counting command completion responses is unreliable, since some of the commands may be skipped and thus not produce a completion message.

56.2.5. Function Call

The Function Call sub-protocol allows the client to request a direct call of any function that exists in the database's `pg_proc` system catalog. The client must have execute permission for the function.

Note

The Function Call sub-protocol is a legacy feature that is probably best avoided in new code. Similar results can be accomplished by setting up a prepared statement that does `SELECT function($1, ...)`. The Function Call cycle can then be replaced with Bind/Execute.

A Function Call cycle is initiated by the frontend sending a FunctionCall message to the backend. The backend then sends one or more response messages depending on the results of the function call, and finally a ReadyForQuery response message. ReadyForQuery informs the frontend that it can safely send a new query or function call.

The possible response messages from the backend are:

ErrorResponse

An error has occurred.

FunctionCallResponse

The function call was completed and returned the result given in the message. (Note that the Function Call protocol can only handle a single scalar result, not a row type or set of results.)

ReadyForQuery

Processing of the function call is complete. ReadyForQuery will always be sent, whether processing terminates successfully or with an error.

NoticeResponse

A warning message has been issued in relation to the function call. Notices are in addition to other responses, i.e., the backend will continue processing the command.

56.2.6. COPY Operations

The `COPY` command allows high-speed bulk data transfer to or from the server. Copy-in and copy-out operations each switch the connection into a distinct sub-protocol, which lasts until the operation is completed.

Copy-in mode (data transfer to the server) is initiated when the backend executes a `COPY FROM STDIN` SQL statement. The backend sends a `CopyInResponse` message to the frontend. The frontend should then send zero or more `CopyData` messages, forming a stream of input data. (The message boundaries are not required to have anything to do with row boundaries, although that is often a reasonable choice.) The frontend can terminate the copy-in mode by sending either a `CopyDone` message (allowing successful termination) or a `CopyFail` message (which will cause the `COPY` SQL statement to fail with an error). The backend then reverts to the command-processing mode it was in before the `COPY` started, which will be either simple or extended query protocol. It will next send either `CommandComplete` (if successful) or `ErrorResponse` (if not).

In the event of a backend-detected error during copy-in mode (including receipt of a `CopyFail` message), the backend will issue an `ErrorResponse` message. If the `COPY` command was issued via an extended-query message, the backend will now discard frontend messages until a `Sync` message is received, then it will issue `ReadyForQuery` and return to normal processing. If the `COPY` command was issued in a simple Query message, the rest of that message is discarded and `ReadyForQuery` is issued. In either case, any subsequent `CopyData`, `CopyDone`, or `CopyFail` messages issued by the frontend will simply be dropped.

The backend will ignore `Flush` and `Sync` messages received during copy-in mode. Receipt of any other non-copy message type constitutes an error that will abort the copy-in state as described above. (The exception for `Flush` and `Sync` is for the convenience of client libraries that always send `Flush` or `Sync` after an `Execute` message, without checking whether the command to be executed is a `COPY FROM STDIN`.)

Copy-out mode (data transfer from the server) is initiated when the backend executes a `COPY TO STDOUT` SQL statement. The backend sends a `CopyOutResponse` message to the frontend, followed by zero or more `CopyData` messages (always one per row), followed by `CopyDone`. The backend then reverts to the command-processing mode it was in before the `COPY` started, and sends `CommandComplete`. The frontend cannot abort the transfer (except by closing the connection or issuing a `Cancel` request), but it can discard unwanted `CopyData` and `CopyDone` messages.

In the event of a backend-detected error during copy-out mode, the backend will issue an `ErrorResponse` message and revert to normal processing. The frontend should treat receipt of `ErrorResponse` as terminating the copy-out mode.

It is possible for NoticeResponse and ParameterStatus messages to be interspersed between CopyData messages; frontends must handle these cases, and should be prepared for other asynchronous message types as well (see [Section 56.2.7](#)). Otherwise, any message type other than CopyData or CopyDone may be treated as terminating copy-out mode.

There is another Copy-related mode called copy-both, which allows high-speed bulk data transfer to *and* from the server. Copy-both mode is initiated when a backend in walsender mode executes a `START_REPLICATION` statement. The backend sends a CopyBothResponse message to the frontend. Both the backend and the frontend may then send CopyData messages until either end sends a CopyDone message. After the client sends a CopyDone message, the connection goes from copy-both mode to copy-out mode, and the client may not send any more CopyData messages. Similarly, when the server sends a CopyDone message, the connection goes into copy-in mode, and the server may not send any more CopyData messages. After both sides have sent a CopyDone message, the copy mode is terminated, and the backend reverts to the command-processing mode. In the event of a backend-detected error during copy-both mode, the backend will issue an ErrorResponse message, discard frontend messages until a Sync message is received, and then issue ReadyForQuery and return to normal processing. The frontend should treat receipt of ErrorResponse as terminating the copy in both directions; no CopyDone should be sent in this case. See [Section 56.4](#) for more information on the subprotocol transmitted over copy-both mode.

The CopyInResponse, CopyOutResponse and CopyBothResponse messages include fields that inform the frontend of the number of columns per row and the format codes being used for each column. (As of the present implementation, all columns in a given `COPY` operation will use the same format, but the message design does not assume this.)

56.2.7. Asynchronous Operations

There are several cases in which the backend will send messages that are not specifically prompted by the frontend's command stream. Frontends must be prepared to deal with these messages at any time, even when not engaged in a query. At minimum, one should check for these cases before beginning to read a query response.

It is possible for NoticeResponse messages to be generated due to outside activity; for example, if the database administrator commands a “fast” database shutdown, the backend will send a NoticeResponse indicating this fact before closing the connection. Accordingly, frontends should always be prepared to accept and display NoticeResponse messages, even when the connection is nominally idle.

ParameterStatus messages will be generated whenever the active value changes for any of the parameters the backend believes the frontend should know about. Most commonly this occurs in response to a `SET SQL` command executed by the frontend, and this case is effectively synchronous — but it is also possible for parameter status changes to occur because the administrator changed a configuration file and then sent the `SIGHUP` signal to the server. Also, if a `SET` command is rolled back, an appropriate ParameterStatus message will be generated to report the current effective value.

At present there is a hard-wired set of parameters for which ParameterStatus will be generated. They are:

<code>application_name</code>	<code>is_superuser</code>
<code>client_encoding</code>	<code>scram_iterations</code>
<code>DateStyle</code>	<code>server_encoding</code>
<code>default_transaction_read_only</code>	<code>server_version</code>
<code>in_hot_standby</code>	<code>session_authorization</code>
<code>integer_datetimes</code>	<code>standard_conforming_strings</code>
<code>IntervalStyle</code>	<code>TimeZone</code>

(`default_transaction_read_only` and `in_hot_standby` were not reported by releases before 14; `scram_iterations` was not reported by releases before 16.) Note that `server_version`, `server_encoding` and `integer_datetimes` are pseudo-parameters that cannot change after startup. This set might change in the future, or even become configurable. Accordingly, a frontend should simply ignore ParameterStatus for parameters that it does not understand or care about.

If a frontend issues a `LISTEN` command, then the backend will send a `NotificationResponse` message (not to be confused with `NoticeResponse`!) whenever a `NOTIFY` command is executed for the same channel name.

Note

At present, `NotificationResponse` can only be sent outside a transaction, and thus it will not occur in the middle of a command-response series, though it might occur just before `ReadyForQuery`. It is unwise to design frontend logic that assumes that, however. Good practice is to be able to accept `NotificationResponse` at any point in the protocol.

56.2.8. Canceling Requests in Progress

During the processing of a query, the frontend might request cancellation of the query. The cancel request is not sent directly on the open connection to the backend for reasons of implementation efficiency: we don't want to have the backend constantly checking for new input from the frontend during query processing. Cancel requests should be relatively infrequent, so we make them slightly cumbersome in order to avoid a penalty in the normal case.

To issue a cancel request, the frontend opens a new connection to the server and sends a `CancelRequest` message, rather than the `StartupMessage` message that would ordinarily be sent across a new connection. The server will process this request and then close the connection. For security reasons, no direct reply is made to the cancel request message.

A `CancelRequest` message will be ignored unless it contains the same key data (PID and secret key) passed to the frontend during connection start-up. If the request matches the PID and secret key for a currently executing backend, the processing of the current query is aborted. (In the existing implementation, this is done by sending a special signal to the backend process that is processing the query.)

The cancellation signal might or might not have any effect — for example, if it arrives after the backend has finished processing the query, then it will have no effect. If the cancellation is effective, it results in the current command being terminated early with an error message.

The upshot of all this is that for reasons of both security and efficiency, the frontend has no direct way to tell whether a cancel request has succeeded. It must continue to wait for the backend to respond to the query. Issuing a cancel simply improves the odds that the current query will finish soon, and improves the odds that it will fail with an error message instead of succeeding.

Since the cancel request is sent across a new connection to the server and not across the regular frontend/backend communication link, it is possible for the cancel request to be issued by any process, not just the frontend whose query is to be canceled. This might provide additional flexibility when building multiple-process applications. It also introduces a security risk, in that unauthorized persons might try to cancel queries. The security risk is addressed by requiring a dynamically generated secret key to be supplied in cancel requests.

56.2.9. Termination

The normal, graceful termination procedure is that the frontend sends a `Terminate` message and immediately closes the connection. On receipt of this message, the backend closes the connection and terminates.

In rare cases (such as an administrator-commanded database shutdown) the backend might disconnect without any frontend request to do so. In such cases the backend will attempt to send an error or notice message giving the reason for the disconnection before it closes the connection.

Other termination scenarios arise from various failure cases, such as core dump at one end or the other, loss of the communications link, loss of message-boundary synchronization, etc. If either frontend or

backend sees an unexpected closure of the connection, it should clean up and terminate. The frontend has the option of launching a new backend by recontacting the server if it doesn't want to terminate itself. Closing the connection is also advisable if an unrecognizable message type is received, since this probably indicates loss of message-boundary sync.

For either normal or abnormal termination, any open transaction is rolled back, not committed. One should note however that if a frontend disconnects while a non-`SELECT` query is being processed, the backend will probably finish the query before noticing the disconnection. If the query is outside any transaction block (`BEGIN ... COMMIT` sequence) then its results might be committed before the disconnection is recognized.

56.2.10. SSL Session Encryption

If Postgres Pro was built with SSL support, frontend/backend communications can be encrypted using SSL. This provides communication security in environments where attackers might be able to capture the session traffic. For more information on encrypting Postgres Pro sessions with SSL, see [Section 18.9](#).

To initiate an SSL-encrypted connection, the frontend initially sends an `SSLRequest` message rather than a `StartupMessage`. The server then responds with a single byte containing `s` or `n`, indicating that it is willing or unwilling to perform SSL, respectively. The frontend might close the connection at this point if it is dissatisfied with the response. To continue after `s`, perform an SSL startup handshake (not described here, part of the SSL specification) with the server. If this is successful, continue with sending the usual `StartupMessage`. In this case the `StartupMessage` and all subsequent data will be SSL-encrypted. To continue after `n`, send the usual `StartupMessage` and proceed without encryption. (Alternatively, it is permissible to issue a `GSSENCRequest` message after an `n` response to try to use GSSAPI encryption instead of SSL.)

The frontend should also be prepared to handle an `ErrorMessage` response to `SSLRequest` from the server. The frontend should not display this error message to the user/application, since the server has not been authenticated ([CVE-2024-10977](#)). In this case the connection must be closed, but the frontend might choose to open a fresh connection and proceed without requesting SSL.

When SSL encryption can be performed, the server is expected to send only the single `s` byte and then wait for the frontend to initiate an SSL handshake. If additional bytes are available to read at this point, it likely means that a man-in-the-middle is attempting to perform a buffer-stuffing attack ([CVE-2021-23222](#)). Frontends should be coded either to read exactly one byte from the socket before turning the socket over to their SSL library, or to treat it as a protocol violation if they find they have read additional bytes.

Likewise the server expects the client to not begin the SSL negotiation until it receives the server's single byte response to the SSL request. If the client begins the SSL negotiation immediately without waiting for the server response to be received it can reduce connection latency by one round-trip. However this comes at the cost of not being able to handle the case where the server sends a negative response to the SSL request. In that case instead of continuing with either GSSAPI or an unencrypted connection or a protocol error the server will simply disconnect.

An initial `SSLRequest` can also be used in a connection that is being opened to send a `CancelRequest` message.

A second alternate way to initiate SSL encryption is available. The server will recognize connections which immediately begin SSL negotiation without any previous `SSLRequest` packets. Once the SSL connection is established the server will expect a normal startup-request packet and continue negotiation over the encrypted channel. In this case any other requests for encryption will be refused. This method is not preferred for general purpose tools as it cannot negotiate the best connection encryption available or handle unencrypted connections. However it is useful for environments where both the server and client are controlled together. In that case it avoids one round trip of latency and allows the use of network tools that depend on standard SSL connections. When using SSL connections in this style the client is required to use the ALPN extension defined by [RFC 7301](#) to protect against protocol confusion attacks. The Postgres Pro protocol is "postgresql" as registered at [IANA TLS ALPN Protocol IDs](#) registry.

While the protocol itself does not provide a way for the server to force SSL encryption, the administrator can configure the server to reject unencrypted sessions as a byproduct of authentication checking.

56.2.11. GSSAPI Session Encryption

If Postgres Pro was built with GSSAPI support, frontend/backend communications can be encrypted using GSSAPI. This provides communication security in environments where attackers might be able to capture the session traffic. For more information on encrypting Postgres Pro sessions with GSSAPI, see [Section 18.10](#).

To initiate a GSSAPI-encrypted connection, the frontend initially sends a `GSSAPIRequest` message rather than a `StartupMessage`. The server then responds with a single byte containing `G` or `N`, indicating that it is willing or unwilling to perform GSSAPI encryption, respectively. The frontend might close the connection at this point if it is dissatisfied with the response. To continue after `G`, using the GSSAPI C bindings as discussed in [RFC 2744](#) or equivalent, perform a GSSAPI initialization by calling `gss_init_sec_context()` in a loop and sending the result to the server, starting with an empty input and then with each result from the server, until it returns no output. When sending the results of `gss_init_sec_context()` to the server, prepend the length of the message as a four byte integer in network byte order. To continue after `N`, send the usual `StartupMessage` and proceed without encryption. (Alternatively, it is permissible to issue an `SSLRequest` message after an `N` response to try to use SSL encryption instead of GSSAPI.)

The frontend should also be prepared to handle an `ErrorMessage` response to `GSSAPIRequest` from the server. The frontend should not display this error message to the user/application, since the server has not been authenticated ([CVE-2024-10977](#)). In this case the connection must be closed, but the frontend might choose to open a fresh connection and proceed without requesting GSSAPI encryption.

When GSSAPI encryption can be performed, the server is expected to send only the single `G` byte and then wait for the frontend to initiate a GSSAPI handshake. If additional bytes are available to read at this point, it likely means that a man-in-the-middle is attempting to perform a buffer-stuffing attack ([CVE-2021-23222](#)). Frontends should be coded either to read exactly one byte from the socket before turning the socket over to their GSSAPI library, or to treat it as a protocol violation if they find they have read additional bytes.

An initial `GSSAPIRequest` can also be used in a connection that is being opened to send a `CancelRequest` message.

Once GSSAPI encryption has been successfully established, use `gss_wrap()` to encrypt the usual `StartupMessage` and all subsequent data, prepending the length of the result from `gss_wrap()` as a four byte integer in network byte order to the actual encrypted payload. Note that the server will only accept encrypted packets from the client which are less than 16kB; `gss_wrap_size_limit()` should be used by the client to determine the size of the unencrypted message which will fit within this limit and larger messages should be broken up into multiple `gss_wrap()` calls. Typical segments are 8kB of unencrypted data, resulting in encrypted packets of slightly larger than 8kB but well within the 16kB maximum. The server can be expected to not send encrypted packets of larger than 16kB to the client.

While the protocol itself does not provide a way for the server to force GSSAPI encryption, the administrator can configure the server to reject unencrypted sessions as a byproduct of authentication checking.

56.3. SASL Authentication

SASL is a framework for authentication in connection-oriented protocols. At the moment, Postgres Pro implements two SASL authentication mechanisms, SCRAM-SHA-256 and SCRAM-SHA-256-PLUS. More might be added in the future. The below steps illustrate how SASL authentication is performed in general, while the next subsection gives more details on SCRAM-SHA-256 and SCRAM-SHA-256-PLUS.

SASL Authentication Message Flow

1. To begin a SASL authentication exchange, the server sends an `AuthenticationSASL` message. It includes a list of SASL authentication mechanisms that the server can accept, in the server's preferred order.

2. The client selects one of the supported mechanisms from the list, and sends a `SASLInitialResponse` message to the server. The message includes the name of the selected mechanism, and an optional Initial Client Response, if the selected mechanism uses that.
3. One or more server-challenge and client-response message will follow. Each server-challenge is sent in an `AuthenticationSASLContinue` message, followed by a response from client in a `SASLResponse` message. The particulars of the messages are mechanism specific.
4. Finally, when the authentication exchange is completed successfully, the server sends an `AuthenticationSASLFinal` message, followed immediately by an `AuthenticationOk` message. The `AuthenticationSASLFinal` contains additional server-to-client data, whose content is particular to the selected authentication mechanism. If the authentication mechanism doesn't use additional data that's sent at completion, the `AuthenticationSASLFinal` message is not sent.

On error, the server can abort the authentication at any stage, and send an `ErrorMessage`.

56.3.1. SCRAM-SHA-256 Authentication

The implemented SASL mechanisms at the moment are `SCRAM-SHA-256` and its variant with channel binding `SCRAM-SHA-256-PLUS`. They are described in detail in [RFC 7677](#) and [RFC 5802](#).

When `SCRAM-SHA-256` is used in Postgres Pro, the server will ignore the user name that the client sends in the `client-first-message`. The user name that was already sent in the startup message is used instead. Postgres Pro supports multiple character encodings, while SCRAM dictates UTF-8 to be used for the user name, so it might be impossible to represent the Postgres Pro user name in UTF-8.

The SCRAM specification dictates that the password is also in UTF-8, and is processed with the *SASLprep* algorithm. Postgres Pro, however, does not require UTF-8 to be used for the password. When a user's password is set, it is processed with SASLprep as if it was in UTF-8, regardless of the actual encoding used. However, if it is not a legal UTF-8 byte sequence, or it contains UTF-8 byte sequences that are prohibited by the SASLprep algorithm, the raw password will be used without SASLprep processing, instead of throwing an error. This allows the password to be normalized when it is in UTF-8, but still allows a non-UTF-8 password to be used, and doesn't require the system to know which encoding the password is in.

Channel binding is supported in Postgres Pro builds with SSL support. The SASL mechanism name for SCRAM with channel binding is `SCRAM-SHA-256-PLUS`. The channel binding type used by Postgres Pro is `tls-server-end-point`.

In SCRAM without channel binding, the server chooses a random number that is transmitted to the client to be mixed with the user-supplied password in the transmitted password hash. While this prevents the password hash from being successfully retransmitted in a later session, it does not prevent a fake server between the real server and client from passing through the server's random value and successfully authenticating.

SCRAM with channel binding prevents such man-in-the-middle attacks by mixing the signature of the server's certificate into the transmitted password hash. While a fake server can retransmit the real server's certificate, it doesn't have access to the private key matching that certificate, and therefore cannot prove it is the owner, causing SSL connection failure.

Example

1. The server sends an `AuthenticationSASL` message. It includes a list of SASL authentication mechanisms that the server can accept. This will be `SCRAM-SHA-256-PLUS` and `SCRAM-SHA-256` if the server is built with SSL support, or else just the latter.
2. The client responds by sending a `SASLInitialResponse` message, which indicates the chosen mechanism, `SCRAM-SHA-256` or `SCRAM-SHA-256-PLUS`. (A client is free to choose either mechanism, but for better security it should choose the channel-binding variant if it can support it.) In the Initial Client response field, the message contains the SCRAM `client-first-message`. The `client-first-message` also contains the channel binding type chosen by the client.

3. Server sends an `AuthenticationSASLContinue` message, with a SCRAM `server-first-message` as the content.
4. Client sends a `SASLResponse` message, with SCRAM `client-final-message` as the content.
5. Server sends an `AuthenticationSASLFinal` message, with the SCRAM `server-final-message`, followed immediately by an `AuthenticationOk` message.

56.4. Streaming Replication Protocol

To initiate streaming replication, the frontend sends the `replication` parameter in the startup message. A Boolean value of `true` (or `on`, `yes`, `1`) tells the backend to go into physical replication walsender mode, wherein a small set of replication commands, shown below, can be issued instead of SQL statements.

Passing `database` as the value for the `replication` parameter instructs the backend to go into logical replication walsender mode, connecting to the database specified in the `dbname` parameter. In logical replication walsender mode, the replication commands shown below as well as normal SQL commands can be issued.

In either physical replication or logical replication walsender mode, only the simple query protocol can be used.

For the purpose of testing replication commands, you can make a replication connection via `psql` or any other libpq-using tool with a connection string including the `replication` option, e.g.:

```
psql "dbname=postgres replication=database" -c "IDENTIFY_SYSTEM;"
```

However, it is often more useful to use [pg_receivewal](#) (for physical replication) or [pg_recvlogical](#) (for logical replication).

Replication commands are logged in the server log when [log_replication_commands](#) is enabled.

The commands accepted in replication mode are:

`IDENTIFY_SYSTEM`

Requests the server to identify itself. Server replies with a result set of a single row, containing four fields:

`systemid` (text)

The unique system identifier identifying the cluster. This can be used to check that the base backup used to initialize the standby came from the same cluster.

`timeline` (int8)

Current timeline ID. Also useful to check that the standby is consistent with the primary.

`xlogpos` (text)

Current WAL flush location. Useful to get a known location in the write-ahead log where streaming can start.

`dbname` (text)

Database connected to or null.

`SHOW name`

Requests the server to send the current setting of a run-time parameter. This is similar to the SQL command [SHOW](#).

`name`

The name of a run-time parameter. Available parameters are documented in [Chapter 19](#).

`TIMELINE_HISTORY tli`

Requests the server to send over the timeline history file for timeline `tli`. Server replies with a result set of a single row, containing two fields. While the fields are labeled as `text`, they effectively return raw bytes, with no encoding conversion:

`filename (text)`

File name of the timeline history file, e.g., `00000002.history`.

`content (text)`

Contents of the timeline history file.

`CREATE_REPLICATION_SLOT slot_name [TEMPORARY] { PHYSICAL | LOGICAL output_plugin } [(option [, ...])]`

Create a physical or logical replication slot. See [Section 26.2.6](#) for more about replication slots.

`slot_name`

The name of the slot to create. Must be a valid replication slot name (see [Section 26.2.6.1](#)).

`output_plugin`

The name of the output plugin used for logical decoding (see [Section 50.6](#)).

`TEMPORARY`

Specify that this replication slot is a temporary one. Temporary slots are not saved to disk and are automatically dropped on error or when the session has finished.

The following options are supported:

`TWO_PHASE [boolean]`

If true, this logical replication slot supports decoding of two-phase commit. With this option, commands related to two-phase commit such as `PREPARE TRANSACTION`, `COMMIT PREPARED` and `ROLLBACK PREPARED` are decoded and transmitted. The transaction will be decoded and transmitted at `PREPARE TRANSACTION` time. The default is false.

`RESERVE_WAL [boolean]`

If true, this physical replication slot reserves WAL immediately. Otherwise, WAL is only reserved upon connection from a streaming replication client. The default is false.

`SNAPSHOT { 'export' | 'use' | 'nothing' }`

Decides what to do with the snapshot created during logical slot initialization. `'export'`, which is the default, will export the snapshot for use in other sessions. This option can't be used inside a transaction. `'use'` will use the snapshot for the current transaction executing the command. This option must be used in a transaction, and `CREATE_REPLICATION_SLOT` must be the first command run in that transaction. Finally, `'nothing'` will just use the snapshot for logical decoding as normal but won't do anything else with it.

`FAILOVER [boolean]`

If true, the slot is enabled to be synced to the standbys so that logical replication can be resumed after failover. The default is false.

In response to this command, the server will send a one-row result set containing the following fields:

`slot_name (text)`

The name of the newly-created replication slot.

`consistent_point (text)`

The WAL location at which the slot became consistent. This is the earliest location from which streaming can start on this replication slot.

`snapshot_name (text)`

The identifier of the snapshot exported by the command. The snapshot is valid until a new command is executed on this connection or the replication connection is closed. Null if the created slot is physical.

`output_plugin (text)`

The name of the output plugin used by the newly-created replication slot. Null if the created slot is physical.

```
CREATE_REPLICATION_SLOT slot_name [ TEMPORARY ] { PHYSICAL [ RESERVE_WAL ] | LOGICAL output_plugin [ EXPORT_SNAPSHOT | NOEXPORT_SNAPSHOT | USE_SNAPSHOT | TWO_PHASE ] }
```

For compatibility with older releases, this alternative syntax for the `CREATE_REPLICATION_SLOT` command is still supported.

```
ALTER_REPLICATION_SLOT slot_name ( option [, ...] )
```

Change the definition of a replication slot. See [Section 26.2.6](#) for more about replication slots. This command is currently only supported for logical replication slots.

slot_name

The name of the slot to alter. Must be a valid replication slot name (see [Section 26.2.6.1](#)).

The following option is supported:

```
FAILOVER [ boolean ]
```

If true, the slot is enabled to be synced to the standbys so that logical replication can be resumed after failover.

```
READ_REPLICATION_SLOT slot_name
```

Read some information associated with a replication slot. Returns a tuple with `NULL` values if the replication slot does not exist. This command is currently only supported for physical replication slots.

In response to this command, the server will return a one-row result set, containing the following fields:

`slot_type (text)`

The replication slot's type, either `physical` or `NULL`.

`restart_lsn (text)`

The replication slot's `restart_lsn`.

`restart_tli (int8)`

The timeline ID associated with `restart_lsn`, following the current timeline history.

```
START_REPLICATION [ SLOT slot_name ] [ PHYSICAL ] XXX/XXX [ TIMELINE tli ]
```

Instructs server to start streaming WAL, starting at WAL location `XXX/XXX`. If `TIMELINE` option is specified, streaming starts on timeline `tli`; otherwise, the server's current timeline is selected. The server can reply with an error, for example if the requested section of WAL has already been recycled.

On success, the server responds with a `CopyBothResponse` message, and then starts to stream WAL to the frontend.

If a slot's name is provided via `slot_name`, it will be updated as replication progresses so that the server knows which WAL segments, and if `hot_standby_feedback` is on which transactions, are still needed by the standby.

If the client requests a timeline that's not the latest but is part of the history of the server, the server will stream all the WAL on that timeline starting from the requested start point up to the point where the server switched to another timeline. If the client requests streaming at exactly the end of an old timeline, the server skips COPY mode entirely.

After streaming all the WAL on a timeline that is not the latest one, the server will end streaming by exiting the COPY mode. When the client acknowledges this by also exiting COPY mode, the server sends a result set with one row and two columns, indicating the next timeline in this server's history. The first column is the next timeline's ID (type `int8`), and the second column is the WAL location where the switch happened (type `text`). Usually, the switch position is the end of the WAL that was streamed, but there are corner cases where the server can send some WAL from the old timeline that it has not itself replayed before promoting. Finally, the server sends two `CommandComplete` messages (one that ends the `CopyData` and the other ends the `START_REPLICATION` itself), and is ready to accept a new command.

WAL data is sent as a series of `CopyData` messages; see [Section 56.6](#) and [Section 56.7](#) for details. (This allows other information to be intermixed; in particular the server can send an `ErrorResponse` message if it encounters a failure after beginning to stream.) The payload of each `CopyData` message from server to the client contains a message of one of the following formats:

XLogData (B)

Byte1('w')

Identifies the message as WAL data.

Int64

The starting point of the WAL data in this message.

Int64

The current end of WAL on the server.

Int64

The server's system clock at the time of transmission, as microseconds since midnight on 2000-01-01.

Byte_n

A section of the WAL data stream.

A single WAL record is never split across two `XLogData` messages. When a WAL record crosses a WAL page boundary, and is therefore already split using continuation records, it can be split at the page boundary. In other words, the first main WAL record and its continuation records can be sent in different `XLogData` messages.

Primary keepalive message (B)

Byte1('k')

Identifies the message as a sender keepalive.

Int64

The current end of WAL on the server.

Int64

The server's system clock at the time of transmission, as microseconds since midnight on 2000-01-01.

Byte1

1 means that the client should reply to this message as soon as possible, to avoid a timeout disconnect. 0 otherwise.

The receiving process can send replies back to the sender at any time, using one of the following message formats (also in the payload of a CopyData message):

Standby status update (F)

Byte1('r')

Identifies the message as a receiver status update.

Int64

The location of the last WAL byte + 1 received and written to disk in the standby.

Int64

The location of the last WAL byte + 1 flushed to disk in the standby.

Int64

The location of the last WAL byte + 1 applied in the standby.

Int64

The client's system clock at the time of transmission, as microseconds since midnight on 2000-01-01.

Byte1

If 1, the client requests the server to reply to this message immediately. This can be used to ping the server, to test if the connection is still healthy.

Hot standby feedback message (F)

Byte1('h')

Identifies the message as a hot standby feedback message.

Int64

The client's system clock at the time of transmission, as microseconds since midnight on 2000-01-01.

Int32

The standby's current global `xmin`, excluding the `catalog_xmin` from any replication slots. If both this value and the following `catalog_xmin` are 0, this is treated as a notification that hot standby feedback will no longer be sent on this connection. Later non-zero messages may reinstate the feedback mechanism.

Int32

The epoch of the global `xmin` xid on the standby.

Int32

The lowest `catalog_xmin` of any replication slots on the standby. Set to 0 if no `catalog_xmin` exists on the standby or if hot standby feedback is being disabled.

Int32

The epoch of the `catalog_xmin xid` on the standby.

```
START_REPLICATION SLOT slot_name LOGICAL XXX/XXX [ ( option_name [ option_value ] [, ...] ) ]
```

Instructs server to start streaming WAL for logical replication, starting at either WAL location `XXX/XXX` or the slot's `confirmed_flush_lsn` (see [Section 55.20](#)), whichever is greater. This behavior makes it easier for clients to avoid updating their local LSN status when there is no data to process. However, starting at a different LSN than requested might not catch certain kinds of client errors; so the client may wish to check that `confirmed_flush_lsn` matches its expectations before issuing `START_REPLICATION`.

The server can reply with an error, for example if the slot does not exist. On success, the server responds with a `CopyBothResponse` message, and then starts to stream WAL to the frontend.

The messages inside the `CopyBothResponse` messages are of the same format documented for `START_REPLICATION ... PHYSICAL`, including two `CommandComplete` messages.

The output plugin associated with the selected slot is used to process the output for streaming.

SLOT slot_name

The name of the slot to stream changes from. This parameter is required, and must correspond to an existing logical replication slot created with `CREATE_REPLICATION_SLOT` in `LOGICAL` mode.

XXX/XXX

The WAL location to begin streaming at.

option_name

The name of an option passed to the slot's logical decoding output plugin. See [Section 56.5](#) for options that are accepted by the standard (`pgoutput`) plugin.

option_value

Optional value, in the form of a string constant, associated with the specified option.

```
DROP_REPLICATION_SLOT slot_name [ WAIT ]
```

Drops a replication slot, freeing any reserved server-side resources.

slot_name

The name of the slot to drop.

WAIT

This option causes the command to wait if the slot is active until it becomes inactive, instead of the default behavior of raising an error.

```
UPLOAD_MANIFEST
```

Uploads a backup manifest in preparation for taking an incremental backup.

```
BASE_BACKUP [ ( option [, ...] ) ]
```

Instructs the server to start streaming a base backup. The system will automatically be put in backup mode before the backup is started, and taken out of it when the backup is complete. The following options are accepted:

`LABEL 'label'`

Sets the label of the backup. If none is specified, a backup label of `base backup` will be used. The quoting rules for the label are the same as a standard SQL string with [standard_conforming_strings](#) turned on.

TARGET *'target'*

Tells the server where to send the backup. If the target is `client`, which is the default, the backup data is sent to the client. If it is `server`, the backup data is written to the server at the pathname specified by the `TARGET_DETAIL` option. If it is `blackhole`, the backup data is not sent anywhere; it is simply discarded.

The `server` target requires superuser privilege or being granted the `pg_write_server_files` role.

TARGET_DETAIL *'detail'*

Provides additional information about the backup target.

Currently, this option can only be used when the backup target is `server`. It specifies the server directory to which the backup should be written.

PROGRESS [*boolean*]

If set to true, request information required to generate a progress report. This will send back an approximate size in the header of each tablespace, which can be used to calculate how far along the stream is done. This is calculated by enumerating all the file sizes once before the transfer is even started, and might as such have a negative impact on the performance. In particular, it might take longer before the first data is streamed. Since the database files can change during the backup, the size is only approximate and might both grow and shrink between the time of approximation and the sending of the actual files. The default is false.

CHECKPOINT { *'fast'* | *'spread'* }

Sets the type of checkpoint to be performed at the beginning of the base backup. The default is `spread`.

WAL [*boolean*]

If set to true, include the necessary WAL segments in the backup. This will include all the files between start and stop backup in the `pg_wal` directory of the base directory tar file. The default is false.

WAIT [*boolean*]

If set to true, the backup will wait until the last required WAL segment has been archived, or emit a warning if WAL archiving is not enabled. If false, the backup will neither wait nor warn, leaving the client responsible for ensuring the required log is available. The default is true.

COMPRESSION *'method'*

Instructs the server to compress the backup using the specified method. Currently, the supported methods are `gzip`, `lz4`, and `zstd`.

COMPRESSION_DETAIL *detail*

Specifies details for the chosen compression method. This should only be used in conjunction with the `COMPRESSION` option. If the value is an integer, it specifies the compression level. Otherwise, it should be a comma-separated list of items, each of the form *keyword* or *keyword=value*. Currently, the supported keywords are `level`, `long` and `workers`.

The `level` keyword sets the compression level. For `gzip` the compression level should be an integer between 1 and 9 (default `Z_DEFAULT_COMPRESSION` or -1), for `lz4` an integer between 1 and 12 (default 0 for fast compression mode), and for `zstd` an integer between `ZSTD_minCLevel()` (usually -131072) and `ZSTD_maxCLevel()` (usually 22), (default `ZSTD_CLEVEL_DEFAULT` or 3).

The `long` keyword enables long-distance matching mode, for improved compression ratio, at the expense of higher memory use. Long-distance mode is supported only for `zstd`.

The `workers` keyword sets the number of threads that should be used for parallel compression. Parallel compression is supported only for `zstd`.

`MAX_RATE` *rate*

Limit (throttle) the maximum amount of data transferred from server to client per unit of time. The expected unit is kilobytes per second. If this option is specified, the value must either be equal to zero or it must fall within the range from 32 kB through 1 GB (inclusive). If zero is passed or the option is not specified, no restriction is imposed on the transfer.

`TABLESPACE_MAP` [*boolean*]

If true, include information about symbolic links present in the directory `pg_tblspc` in a file named `tablespace_map`. The tablespace map file includes each symbolic link name as it exists in the directory `pg_tblspc/` and the full path of that symbolic link. The default is false.

`VERIFY_CHECKSUMS` [*boolean*]

If true, checksums are verified during a base backup if they are enabled. If false, this is skipped. The default is true.

`MANIFEST` *manifest_option*

When this option is specified with a value of `yes` or `force-encode`, a backup manifest is created and sent along with the backup. The manifest is a list of every file present in the backup with the exception of any WAL files that may be included. It also stores the size, last modification time, and optionally a checksum for each file. A value of `force-encode` forces all filenames to be hex-encoded; otherwise, this type of encoding is performed only for files whose names are non-UTF8 octet sequences. `force-encode` is intended primarily for testing purposes, to be sure that clients which read the backup manifest can handle this case. For compatibility with previous releases, the default is `MANIFEST 'no'`.

`MANIFEST_CHECKSUMS` *checksum_algorithm*

Specifies the checksum algorithm that should be applied to each file included in the backup manifest. Currently, the available algorithms are `NONE`, `CRC32C`, `SHA224`, `SHA256`, `SHA384`, and `SHA512`. The default is `CRC32C`.

`INCREMENTAL`

Requests an incremental backup. The `UPLOAD_MANIFEST` command must be executed before running a base backup with this option.

When the backup is started, the server will first send two ordinary result sets, followed by one or more `CopyOutResponse` results.

The first ordinary result set contains the starting position of the backup, in a single row with two columns. The first column contains the start position given in `XLogRecPtr` format, and the second column contains the corresponding timeline ID.

The second ordinary result set has one row for each tablespace. The fields in this row are:

`spcoid` (*oid*)

The OID of the tablespace, or null if it's the base directory.

`spclocation` (*text*)

The full path of the tablespace directory, or null if it's the base directory.

`size` (*int8*)

The approximate size of the tablespace, in kilobytes (1024 bytes), if progress report has been requested; otherwise it's null.

After the second regular result set, a `CopyOutResponse` will be sent. The payload of each `CopyData` message will contain a message in one of the following formats:

new archive (B)

Byte1('n')

Identifies the message as indicating the start of a new archive. There will be one archive for the main data directory and one for each additional tablespace; each will use tar format (following the “ustar interchange format” specified in the POSIX 1003.1-2008 standard).

String

The file name for this archive.

String

For the main data directory, an empty string. For other tablespaces, the full path to the directory from which this archive was created.

manifest (B)

Byte1('m')

Identifies the message as indicating the start of the backup manifest.

archive or manifest data (B)

Byte1('d')

Identifies the message as containing archive or manifest data.

Byte_n

Data bytes.

progress report (B)

Byte1('p')

Identifies the message as a progress report.

Int64

The number of bytes from the current tablespace for which processing has been completed.

After the `CopyOutResponse`, or all such responses, have been sent, a final ordinary result set will be sent, containing the WAL end position of the backup, in the same format as the start position.

The tar archive for the data directory and each tablespace will contain all files in the directories, regardless of whether they are Postgres Pro files or other files added to the same directory. The only excluded files are:

- `postmaster.pid`
- `postmaster.opts`
- `pg_internal.init` (found in multiple directories)
- Various temporary files and directories created during the operation of the Postgres Pro server, such as any file or directory beginning with `pgsql_tmp` and temporary relations.
- Unlogged relations, except for the init fork which is required to recreate the (empty) unlogged relation on recovery.
- `pg_wal`, including subdirectories. If the backup is run with WAL files included, a synthesized version of `pg_wal` will be included, but it will only contain the files necessary for the backup to work, not the rest of the contents.
- `pg_dynshmem`, `pg_notify`, `pg_replslot`, `pg_serial`, `pg_snapshots`, `pg_stat_tmp`, and `pg_subtrans` are copied as empty directories (even if they are symbolic links).

- Files other than regular files and directories, such as symbolic links (other than for the directories listed above) and special device and operating system files, are skipped. (Symbolic links in `pg_tblspc` are maintained.)

Owner, group, and file mode are set if the underlying file system on the server supports it.

In all the above commands, when specifying a parameter of type `boolean` the *value* part can be omitted, which is equivalent to specifying `TRUE`.

56.5. Logical Streaming Replication Protocol

This section describes the logical replication protocol, which is the message flow started by the `START_REPLICATION SLOT slot_name LOGICAL` replication command.

The logical streaming replication protocol builds on the primitives of the physical streaming replication protocol.

Postgres Pro logical decoding supports output plugins. `pgoutput` is the standard one used for the built-in logical replication.

56.5.1. Logical Streaming Replication Parameters

Using the `START_REPLICATION` command, `pgoutput` accepts the following options:

`proto_version`

Protocol version. Currently versions 1, 2, 3, and 4 are supported. A valid version is required.

Version 2 is supported only for server version 14 and above, and it allows streaming of large in-progress transactions.

Version 3 is supported only for server version 15 and above, and it allows streaming of two-phase commits.

Version 4 is supported only for server version 16 and above, and it allows streams of large in-progress transactions to be applied in parallel.

`publication_names`

Comma separated list of publication names for which to subscribe (receive changes). The individual publication names are treated as standard objects names and can be quoted the same as needed. At least one publication name is required.

`binary`

Boolean option to use binary transfer mode. Binary mode is faster than the text mode but slightly less robust.

`messages`

Boolean option to enable sending the messages that are written by `pg_logical_emit_message`.

`streaming`

Boolean option to enable streaming of in-progress transactions. It accepts an additional value "parallel" to enable sending extra information with some messages to be used for parallelisation. Minimum protocol version 2 is required to turn it on. Minimum protocol version 4 is required for the "parallel" option.

`two_phase`

Boolean option to enable two-phase transactions. Minimum protocol version 3 is required to turn it on.

origin

Option to send changes by their origin. Possible values are "none" to only send the changes that have no origin associated, or "any" to send the changes regardless of their origin. This can be used to avoid loops (infinite replication of the same data) among replication nodes.

56.5.2. Logical Replication Protocol Messages

The individual protocol messages are discussed in the following subsections. Individual messages are described in [Section 56.9](#).

All top-level protocol messages begin with a message type byte. While represented in code as a character, this is a signed byte with no associated encoding.

Since the streaming replication protocol supplies a message length there is no need for top-level protocol messages to embed a length in their header.

56.5.3. Logical Replication Protocol Message Flow

With the exception of the `START_REPLICATION` command and the replay progress messages, all information flows only from the backend to the frontend.

The logical replication protocol sends individual transactions one by one. This means that all messages between a pair of Begin and Commit messages belong to the same transaction. Similarly, all messages between a pair of Begin Prepare and Prepare messages belong to the same transaction. It also sends changes of large in-progress transactions between a pair of Stream Start and Stream Stop messages. The last stream of such a transaction contains a Stream Commit or Stream Abort message.

Every sent transaction contains zero or more DML messages (Insert, Update, Delete). In case of a cascaded setup it can also contain Origin messages. The origin message indicates that the transaction originated on different replication node. Since a replication node in the scope of logical replication protocol can be pretty much anything, the only identifier is the origin name. It's downstream's responsibility to handle this as needed (if needed). The Origin message is always sent before any DML messages in the transaction.

Every DML message contains a relation OID, identifying the publisher's relation that was acted on. Before the first DML message for a given relation OID, a Relation message will be sent, describing the schema of that relation. Subsequently, a new Relation message will be sent if the relation's definition has changed since the last Relation message was sent for it. (The protocol assumes that the client is capable of remembering this metadata for as many relations as needed.)

Relation messages identify column types by their OIDs. In the case of a built-in type, it is assumed that the client can look up that type OID locally, so no additional data is needed. For a non-built-in type OID, a Type message will be sent before the Relation message, to provide the type name associated with that OID. Thus, a client that needs to specifically identify the types of relation columns should cache the contents of Type messages, and first consult that cache to see if the type OID is defined there. If not, look up the type OID locally.

56.6. Message Data Types

This section describes the base data types used in messages.

$\text{Int}_n(i)$

An n -bit integer in network byte order (most significant byte first). If i is specified it is the exact value that will appear, otherwise the value is variable. Eg. `Int16`, `Int32(42)`.

$\text{Int}_n[k]$

An array of k n -bit integers, each in network byte order. The array length k is always determined by an earlier field in the message. Eg. `Int16[M]`.

String(*s*)

A null-terminated string (C-style string). There is no specific length limitation on strings. If *s* is specified it is the exact value that will appear, otherwise the value is variable. Eg. String, String("user").

Note

There is no predefined limit on the length of a string that can be returned by the backend. Good coding strategy for a frontend is to use an expandable buffer so that anything that fits in memory can be accepted. If that's not feasible, read the full string and discard trailing characters that don't fit into your fixed-size buffer.

Byte_{*n*}(*c*)

Exactly *n* bytes. If the field width *n* is not a constant, it is always determinable from an earlier field in the message. If *c* is specified it is the exact value. Eg. Byte2, Byte1('\n').

56.7. Message Formats

This section describes the detailed format of each message. Each is marked to indicate that it can be sent by a frontend (F), a backend (B), or both (F & B). Notice that although each message includes a byte count at the beginning, the message format is defined so that the message end can be found without reference to the byte count. This aids validity checking. (The CopyData message is an exception, because it forms part of a data stream; the contents of any individual CopyData message cannot be interpretable on their own.)

AuthenticationOk (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(8)

Length of message contents in bytes, including self.

Int32(0)

Specifies that the authentication was successful.

AuthenticationKerberosV5 (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(8)

Length of message contents in bytes, including self.

Int32(2)

Specifies that Kerberos V5 authentication is required.

AuthenticationCleartextPassword (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(8)

Length of message contents in bytes, including self.

Int32(3)

Specifies that a clear-text password is required.

AuthenticationMD5Password (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(12)

Length of message contents in bytes, including self.

Int32(5)

Specifies that an MD5-encrypted password is required.

Byte4

The salt to use when encrypting the password.

AuthenticationGSS (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(8)

Length of message contents in bytes, including self.

Int32(7)

Specifies that GSSAPI authentication is required.

AuthenticationGSSContinue (B)

Byte1('R')

Identifies the message as an authentication request.

Int32

Length of message contents in bytes, including self.

Int32(8)

Specifies that this message contains GSSAPI or SSPI data.

Byte_n

GSSAPI or SSPI authentication data.

AuthenticationSSPI (B)

Byte1('R')

Identifies the message as an authentication request.

Int32(8)

Length of message contents in bytes, including self.

Int32(9)

Specifies that SSPI authentication is required.

AuthenticationSASL (B)

Byte1('R')

Identifies the message as an authentication request.

Int32

Length of message contents in bytes, including self.

Int32(10)

Specifies that SASL authentication is required.

The message body is a list of SASL authentication mechanisms, in the server's order of preference. A zero byte is required as terminator after the last authentication mechanism name. For each mechanism, there is the following:

String

Name of a SASL authentication mechanism.

AuthenticationSASLContinue (B)

Byte1('R')

Identifies the message as an authentication request.

Int32

Length of message contents in bytes, including self.

Int32(11)

Specifies that this message contains a SASL challenge.

Byte_n

SASL data, specific to the SASL mechanism being used.

AuthenticationSASLFinal (B)

Byte1('R')

Identifies the message as an authentication request.

Int32

Length of message contents in bytes, including self.

Int32(12)

Specifies that SASL authentication has completed.

Byte_n

SASL outcome "additional data", specific to the SASL mechanism being used.

BackendKeyData (B)

Byte1('K')

Identifies the message as cancellation key data. The frontend must save these values if it wishes to be able to issue CancelRequest messages later.

Int32(12)

Length of message contents in bytes, including self.

Int32

The process ID of this backend.

Int32

The secret key of this backend.

Bind (F)**Byte1('B')**

Identifies the message as a Bind command.

Int32

Length of message contents in bytes, including self.

String

The name of the destination portal (an empty string selects the unnamed portal).

String

The name of the source prepared statement (an empty string selects the unnamed prepared statement).

Int16

The number of parameter format codes that follow (denoted *c* below). This can be zero to indicate that there are no parameters or that the parameters all use the default format (text); or one, in which case the specified format code is applied to all parameters; or it can equal the actual number of parameters.

Int16[*c*]

The parameter format codes. Each must presently be zero (text) or one (binary).

Int16

The number of parameter values that follow (possibly zero). This must match the number of parameters needed by the query.

Next, the following pair of fields appear for each parameter:

Int32

The length of the parameter value, in bytes (this count does not include itself). Can be zero. As a special case, -1 indicates a NULL parameter value. No value bytes follow in the NULL case.

Byte_{*n*}

The value of the parameter, in the format indicated by the associated format code. *n* is the above length.

After the last parameter, the following fields appear:

Int16

The number of result-column format codes that follow (denoted *R* below). This can be zero to indicate that there are no result columns or that the result columns should all use the default format (text); or one, in which case the specified format code is applied to all result columns (if any); or it can equal the actual number of result columns of the query.

Int16[*R*]

The result-column format codes. Each must presently be zero (text) or one (binary).

BindComplete (B)**Byte1('2')**

Identifies the message as a Bind-complete indicator.

Int32(4)

Length of message contents in bytes, including self.

CancelRequest (F)

Int32(16)

Length of message contents in bytes, including self.

Int32(80877102)

The cancel request code. The value is chosen to contain 1234 in the most significant 16 bits, and 5678 in the least significant 16 bits. (To avoid confusion, this code must not be the same as any protocol version number.)

Int32

The process ID of the target backend.

Int32

The secret key for the target backend.

Close (F)

Byte1('C')

Identifies the message as a Close command.

Int32

Length of message contents in bytes, including self.

Byte1

's' to close a prepared statement; or 'P' to close a portal.

String

The name of the prepared statement or portal to close (an empty string selects the unnamed prepared statement or portal).

CloseComplete (B)

Byte1('3')

Identifies the message as a Close-complete indicator.

Int32(4)

Length of message contents in bytes, including self.

CommandComplete (B)

Byte1('C')

Identifies the message as a command-completed response.

Int32

Length of message contents in bytes, including self.

String

The command tag. This is usually a single word that identifies which SQL command was completed.

For an `INSERT` command, the tag is `INSERT oid rows`, where `rows` is the number of rows inserted. `oid` used to be the object ID of the inserted row if `rows` was 1 and the target table had OIDs, but OIDs system columns are not supported anymore; therefore `oid` is always 0.

For a `DELETE` command, the tag is `DELETE rows` where `rows` is the number of rows deleted.

For an `UPDATE` command, the tag is `UPDATE rows` where *rows* is the number of rows updated.

For a `MERGE` command, the tag is `MERGE rows` where *rows* is the number of rows inserted, updated, or deleted.

For a `SELECT` or `CREATE TABLE AS` command, the tag is `SELECT rows` where *rows* is the number of rows retrieved.

For a `MOVE` command, the tag is `MOVE rows` where *rows* is the number of rows the cursor's position has been changed by.

For a `FETCH` command, the tag is `FETCH rows` where *rows* is the number of rows that have been retrieved from the cursor.

For a `COPY` command, the tag is `COPY rows` where *rows* is the number of rows copied. (Note: the row count appears only in PostgreSQL 8.2 and later.)

CopyData (F & B)

Byte1('d')

Identifies the message as `COPY` data.

Int32

Length of message contents in bytes, including self.

Byte_n

Data that forms part of a `COPY` data stream. Messages sent from the backend will always correspond to single data rows, but messages sent by frontends might divide the data stream arbitrarily.

CopyDone (F & B)

Byte1('c')

Identifies the message as a `COPY`-complete indicator.

Int32(4)

Length of message contents in bytes, including self.

CopyFail (F)

Byte1('f')

Identifies the message as a `COPY`-failure indicator.

Int32

Length of message contents in bytes, including self.

String

An error message to report as the cause of failure.

CopyInResponse (B)

Byte1('G')

Identifies the message as a Start Copy In response. The frontend must now send copy-in data (if not prepared to do so, send a `CopyFail` message).

Int32

Length of message contents in bytes, including self.

Int8

0 indicates the overall `COPY` format is textual (rows separated by newlines, columns separated by separator characters, etc.). 1 indicates the overall copy format is binary (similar to `DataRow` format). See [COPY](#) for more information.

Int16

The number of columns in the data to be copied (denoted N below).

Int16[N]

The format codes to be used for each column. Each must presently be zero (text) or one (binary). All must be zero if the overall copy format is textual.

CopyOutResponse (B)

Byte1('H')

Identifies the message as a Start Copy Out response. This message will be followed by copy-out data.

Int32

Length of message contents in bytes, including self.

Int8

0 indicates the overall `COPY` format is textual (rows separated by newlines, columns separated by separator characters, etc.). 1 indicates the overall copy format is binary (similar to `DataRow` format). See [COPY](#) for more information.

Int16

The number of columns in the data to be copied (denoted N below).

Int16[N]

The format codes to be used for each column. Each must presently be zero (text) or one (binary). All must be zero if the overall copy format is textual.

CopyBothResponse (B)

Byte1('W')

Identifies the message as a Start Copy Both response. This message is used only for Streaming Replication.

Int32

Length of message contents in bytes, including self.

Int8

0 indicates the overall `COPY` format is textual (rows separated by newlines, columns separated by separator characters, etc.). 1 indicates the overall copy format is binary (similar to `DataRow` format). See [COPY](#) for more information.

Int16

The number of columns in the data to be copied (denoted N below).

Int16[N]

The format codes to be used for each column. Each must presently be zero (text) or one (binary). All must be zero if the overall copy format is textual.

DataRow (B)**Byte1('D')**

Identifies the message as a data row.

Int32

Length of message contents in bytes, including self.

Int16

The number of column values that follow (possibly zero).

Next, the following pair of fields appear for each column:

Int32

The length of the column value, in bytes (this count does not include itself). Can be zero. As a special case, -1 indicates a NULL column value. No value bytes follow in the NULL case.

Byte_{*n*}The value of the column, in the format indicated by the associated format code. *n* is the above length.**Describe (F)****Byte1('D')**

Identifies the message as a Describe command.

Int32

Length of message contents in bytes, including self.

Byte1

'S' to describe a prepared statement; or 'P' to describe a portal.

String

The name of the prepared statement or portal to describe (an empty string selects the unnamed prepared statement or portal).

EmptyQueryResponse (B)**Byte1('I')**

Identifies the message as a response to an empty query string. (This substitutes for Command-Complete.)

Int32(4)

Length of message contents in bytes, including self.

ErrorResponse (B)**Byte1('E')**

Identifies the message as an error.

Int32

Length of message contents in bytes, including self.

The message body consists of one or more identified fields, followed by a zero byte as a terminator. Fields can appear in any order. For each field there is the following:

Byte1

A code identifying the field type; if zero, this is the message terminator and no string follows. The presently defined field types are listed in [Section 56.8](#). Since more field types might be added in future, frontends should silently ignore fields of unrecognized type.

String

The field value.

Execute (F)**Byte1('E')**

Identifies the message as an Execute command.

Int32

Length of message contents in bytes, including self.

String

The name of the portal to execute (an empty string selects the unnamed portal).

Int32

Maximum number of rows to return, if portal contains a query that returns rows (ignored otherwise). Zero denotes “no limit”.

Flush (F)**Byte1('H')**

Identifies the message as a Flush command.

Int32(4)

Length of message contents in bytes, including self.

FunctionCall (F)**Byte1('F')**

Identifies the message as a function call.

Int32

Length of message contents in bytes, including self.

Int32

Specifies the object ID of the function to call.

Int16

The number of argument format codes that follow (denoted *c* below). This can be zero to indicate that there are no arguments or that the arguments all use the default format (text); or one, in which case the specified format code is applied to all arguments; or it can equal the actual number of arguments.

Int16[*c*]

The argument format codes. Each must presently be zero (text) or one (binary).

Int16

Specifies the number of arguments being supplied to the function.

Next, the following pair of fields appear for each argument:

Int32

The length of the argument value, in bytes (this count does not include itself). Can be zero. As a special case, -1 indicates a NULL argument value. No value bytes follow in the NULL case.

Byte_n

The value of the argument, in the format indicated by the associated format code. *n* is the above length.

After the last argument, the following field appears:

Int16

The format code for the function result. Must presently be zero (text) or one (binary).

FunctionCallResponse (B)**Byte1('V')**

Identifies the message as a function call result.

Int32

Length of message contents in bytes, including self.

Int32

The length of the function result value, in bytes (this count does not include itself). Can be zero. As a special case, -1 indicates a NULL function result. No value bytes follow in the NULL case.

Byte_n

The value of the function result, in the format indicated by the associated format code. *n* is the above length.

GSSENCRequest (F)**Int32(8)**

Length of message contents in bytes, including self.

Int32(80877104)

The GSSAPI Encryption request code. The value is chosen to contain 1234 in the most significant 16 bits, and 5680 in the least significant 16 bits. (To avoid confusion, this code must not be the same as any protocol version number.)

GSSResponse (F)**Byte1('p')**

Identifies the message as a GSSAPI or SSPI response. Note that this is also used for SASL and password response messages. The exact message type can be deduced from the context.

Int32

Length of message contents in bytes, including self.

Byte_n

GSSAPI/SSPI specific message data.

NegotiateProtocolVersion (B)**Byte1('v')**

Identifies the message as a protocol version negotiation message.

Int32

Length of message contents in bytes, including self.

Int32

Newest minor protocol version supported by the server for the major protocol version requested by the client.

Int32

Number of protocol options not recognized by the server.

Then, for protocol option not recognized by the server, there is the following:

String

The option name.

NoData (B)

Byte1('n')

Identifies the message as a no-data indicator.

Int32(4)

Length of message contents in bytes, including self.

NoticeResponse (B)

Byte1('N')

Identifies the message as a notice.

Int32

Length of message contents in bytes, including self.

The message body consists of one or more identified fields, followed by a zero byte as a terminator. Fields can appear in any order. For each field there is the following:

Byte1

A code identifying the field type; if zero, this is the message terminator and no string follows. The presently defined field types are listed in [Section 56.8](#). Since more field types might be added in future, frontends should silently ignore fields of unrecognized type.

String

The field value.

NotificationResponse (B)

Byte1('A')

Identifies the message as a notification response.

Int32

Length of message contents in bytes, including self.

Int32

The process ID of the notifying backend process.

String

The name of the channel that the notify has been raised on.

String

The “payload” string passed from the notifying process.

ParameterDescription (B)

Byte1('t')

Identifies the message as a parameter description.

Int32

Length of message contents in bytes, including self.

Int16

The number of parameters used by the statement (can be zero).

Then, for each parameter, there is the following:

Int32

Specifies the object ID of the parameter data type.

ParameterStatus (B)

Byte1('S')

Identifies the message as a run-time parameter status report.

Int32

Length of message contents in bytes, including self.

String

The name of the run-time parameter being reported.

String

The current value of the parameter.

Parse (F)

Byte1('P')

Identifies the message as a Parse command.

Int32

Length of message contents in bytes, including self.

String

The name of the destination prepared statement (an empty string selects the unnamed prepared statement).

String

The query string to be parsed.

Int16

The number of parameter data types specified (can be zero). Note that this is not an indication of the number of parameters that might appear in the query string, only the number that the frontend wants to prespecify types for.

Then, for each parameter, there is the following:

Int32

Specifies the object ID of the parameter data type. Placing a zero here is equivalent to leaving the type unspecified.

ParseComplete (B)

Byte1('1')

Identifies the message as a Parse-complete indicator.

Int32(4)

Length of message contents in bytes, including self.

PasswordMessage (F)

Byte1('p')

Identifies the message as a password response. Note that this is also used for GSSAPI, SSPI and SASL response messages. The exact message type can be deduced from the context.

Int32

Length of message contents in bytes, including self.

String

The password (encrypted, if requested).

PortalSuspended (B)

Byte1('s')

Identifies the message as a portal-suspended indicator. Note this only appears if an Execute message's row-count limit was reached.

Int32(4)

Length of message contents in bytes, including self.

Query (F)

Byte1('Q')

Identifies the message as a simple query.

Int32

Length of message contents in bytes, including self.

String

The query string itself.

ReadyForQuery (B)

Byte1('Z')

Identifies the message type. ReadyForQuery is sent whenever the backend is ready for a new query cycle.

Int32(5)

Length of message contents in bytes, including self.

Byte1

Current backend transaction status indicator. Possible values are 'I' if idle (not in a transaction block); 'T' if in a transaction block; or 'E' if in a failed transaction block (queries will be rejected until block is ended).

RowDescription (B)

Byte1('T')

Identifies the message as a row description.

Int32

Length of message contents in bytes, including self.

Int16

Specifies the number of fields in a row (can be zero).

Then, for each field, there is the following:

String

The field name.

Int32

If the field can be identified as a column of a specific table, the object ID of the table; otherwise zero.

Int16

If the field can be identified as a column of a specific table, the attribute number of the column; otherwise zero.

Int32

The object ID of the field's data type.

Int16

The data type size (see `pg_type.typelen`). Note that negative values denote variable-width types.

Int32

The type modifier (see `pg_attribute.atttypmod`). The meaning of the modifier is type-specific.

Int16

The format code being used for the field. Currently will be zero (text) or one (binary). In a Row-Description returned from the statement variant of Describe, the format code is not yet known and will always be zero.

SASLInitialResponse (F)

Byte1('p')

Identifies the message as an initial SASL response. Note that this is also used for GSSAPI, SSPI and password response messages. The exact message type is deduced from the context.

Int32

Length of message contents in bytes, including self.

String

Name of the SASL authentication mechanism that the client selected.

Int32

Length of SASL mechanism specific "Initial Client Response" that follows, or -1 if there is no Initial Response.

Byte_n

SASL mechanism specific "Initial Response".

SASLResponse (F)**Byte1('p')**

Identifies the message as a SASL response. Note that this is also used for GSSAPI, SSPI and password response messages. The exact message type can be deduced from the context.

Int32

Length of message contents in bytes, including self.

Byte_n

SASL mechanism specific message data.

SSLRequest (F)**Int32(8)**

Length of message contents in bytes, including self.

Int32(80877103)

The SSL request code. The value is chosen to contain 1234 in the most significant 16 bits, and 5679 in the least significant 16 bits. (To avoid confusion, this code must not be the same as any protocol version number.)

StartupMessage (F)**Int32**

Length of message contents in bytes, including self.

Int32(196608)

The protocol version number. The most significant 16 bits are the major version number (3 for the protocol described here). The least significant 16 bits are the minor version number (0 for the protocol described here).

The protocol version number is followed by one or more pairs of parameter name and value strings. A zero byte is required as a terminator after the last name/value pair. Parameters can appear in any order. `user` is required, others are optional. Each parameter is specified as:

String

The parameter name. Currently recognized names are:

`user`

The database user name to connect as. Required; there is no default.

`database`

The database to connect to. Defaults to the user name.

`options`

Command-line arguments for the backend. (This is deprecated in favor of setting individual run-time parameters.) Spaces within this string are considered to separate arguments, unless escaped with a backslash (`\`); write `\\` to represent a literal backslash.

`replication`

Used to connect in streaming replication mode, where a small set of replication commands can be issued instead of SQL statements. Value can be `true`, `false`, or `database`, and the default is `false`. See [Section 56.4](#) for details.

In addition to the above, other parameters may be listed. Parameter names beginning with `_pq_` are reserved for use as protocol extensions, while others are treated as run-time parameters to

be set at backend start time. Such settings will be applied during backend start (after parsing the command-line arguments if any) and will act as session defaults.

String

The parameter value.

Sync (F)

Byte1('S')

Identifies the message as a Sync command.

Int32(4)

Length of message contents in bytes, including self.

Terminate (F)

Byte1('X')

Identifies the message as a termination.

Int32(4)

Length of message contents in bytes, including self.

56.8. Error and Notice Message Fields

This section describes the fields that can appear in ErrorResponse and NoticeResponse messages. Each field type has a single-byte identification token. Note that any given field type should appear at most once per message.

S

Severity: the field contents are `ERROR`, `FATAL`, or `PANIC` (in an error message), or `WARNING`, `NOTICE`, `DEBUG`, `INFO`, or `LOG` (in a notice message), or a localized translation of one of these. Always present.

V

Severity: the field contents are `ERROR`, `FATAL`, or `PANIC` (in an error message), or `WARNING`, `NOTICE`, `DEBUG`, `INFO`, or `LOG` (in a notice message). This is identical to the S field except that the contents are never localized. This is present only in messages generated by Postgres Pro versions 9.6 and later.

C

Code: the SQLSTATE code for the error (see [Appendix A](#)). Not localizable. Always present.

M

Message: the primary human-readable error message. This should be accurate but terse (typically one line). Always present.

D

Detail: an optional secondary error message carrying more detail about the problem. Might run to multiple lines.

H

Hint: an optional suggestion what to do about the problem. This is intended to differ from Detail in that it offers advice (potentially inappropriate) rather than hard facts. Might run to multiple lines.

P

Position: the field value is a decimal ASCII integer, indicating an error cursor position as an index into the original query string. The first character has index 1, and positions are measured in characters not bytes.

P

Internal position: this is defined the same as the P field, but it is used when the cursor position refers to an internally generated command rather than the one submitted by the client. The q field will always appear when this field appears.

Q

Internal query: the text of a failed internally-generated command. This could be, for example, an SQL query issued by a PL/pgSQL function.

W

Where: an indication of the context in which the error occurred. Presently this includes a call stack traceback of active procedural language functions and internally-generated queries. The trace is one entry per line, most recent first.

S

Schema name: if the error was associated with a specific database object, the name of the schema containing that object, if any.

T

Table name: if the error was associated with a specific table, the name of the table. (Refer to the schema name field for the name of the table's schema.)

C

Column name: if the error was associated with a specific table column, the name of the column. (Refer to the schema and table name fields to identify the table.)

D

Data type name: if the error was associated with a specific data type, the name of the data type. (Refer to the schema name field for the name of the data type's schema.)

N

Constraint name: if the error was associated with a specific constraint, the name of the constraint. Refer to fields listed above for the associated table or domain. (For this purpose, indexes are treated as constraints, even if they weren't created with constraint syntax.)

F

File: the file name of the source-code location where the error was reported.

L

Line: the line number of the source-code location where the error was reported.

R

Routine: the name of the source-code routine reporting the error.

Note

The fields for schema name, table name, column name, data type name, and constraint name are supplied only for a limited number of error types; see [Appendix A](#). Frontends should not assume that the presence of any of these fields guarantees the presence of another field. Core error sources observe the interrelationships noted above, but user-defined functions may use these fields in other ways. In the same vein, clients should not assume that these fields denote contemporary objects in the current database.

The client is responsible for formatting displayed information to meet its needs; in particular it should break long lines as needed. Newline characters appearing in the error message fields should be treated as paragraph breaks, not line breaks.

56.9. Logical Replication Message Formats

This section describes the detailed format of each logical replication message. These messages are either returned by the replication slot SQL interface or are sent by a walsender. In the case of a walsender, they are encapsulated inside replication protocol WAL messages as described in [Section 56.4](#), and generally obey the same message flow as physical replication.

Begin

Byte1('B')

Identifies the message as a begin message.

Int64 (XLogRecPtr)

The final LSN of the transaction.

Int64 (TimestampTz)

Commit timestamp of the transaction. The value is in number of microseconds since Postgres Pro epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

Message

Byte1('M')

Identifies the message as a logical decoding message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int8

Flags; Either 0 for no flags or 1 if the logical decoding message is transactional.

Int64 (XLogRecPtr)

The LSN of the logical decoding message.

String

The prefix of the logical decoding message.

Int32

Length of the content.

Byte_n

The content of the logical decoding message.

Commit

Byte1('C')

Identifies the message as a commit message.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The LSN of the commit.

Int64 (XLogRecPtr)

The end LSN of the transaction.

Int64 (TimestampTz)

Commit timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Origin

Byte1('O')

Identifies the message as an origin message.

Int64 (XLogRecPtr)

The LSN of the commit on the origin server.

String

Name of the origin.

Note that there can be multiple Origin messages inside a single transaction.

Relation

Byte1('R')

Identifies the message as a relation message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32 (Oid)

OID of the relation.

String

Namespace (empty string for `pg_catalog`).

String

Relation name.

Int8

Replica identity setting for the relation (same as `relreplident` in `pg_class`).

Int16

Number of columns.

Next, the following message part appears for each column included in the publication (except generated columns):

Int8

Flags for the column. Currently can be either 0 for no flags or 1 which marks the column as part of the key.

String

Name of the column.

Int32 (Oid)

OID of the column's data type.

Int32

Type modifier of the column (`atttypmod`).

Type

Byte1('Y')

Identifies the message as a type message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32 (Oid)

OID of the data type.

String

Namespace (empty string for `pg_catalog`).

String

Name of the data type.

Insert

Byte1('I')

Identifies the message as an insert message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32 (Oid)

OID of the relation corresponding to the ID in the relation message.

Byte1('N')

Identifies the following `TupleData` message as a new tuple.

TupleData

`TupleData` message part representing the contents of new tuple.

Update

Byte1('U')

Identifies the message as an update message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32 (Oid)

OID of the relation corresponding to the ID in the relation message.

Byte1('K')

Identifies the following TupleData submessage as a key. This field is optional and is only present if the update changed data in any of the column(s) that are part of the REPLICA IDENTITY index.

Byte1('O')

Identifies the following TupleData submessage as an old tuple. This field is optional and is only present if table in which the update happened has REPLICA IDENTITY set to FULL.

TupleData

TupleData message part representing the contents of the old tuple or primary key. Only present if the previous 'O' or 'K' part is present.

Byte1('N')

Identifies the following TupleData message as a new tuple.

TupleData

TupleData message part representing the contents of a new tuple.

The Update message may contain either a 'K' message part or an 'O' message part or neither of them, but never both of them.

Delete**Byte1('D')**

Identifies the message as a delete message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32 (Oid)

OID of the relation corresponding to the ID in the relation message.

Byte1('K')

Identifies the following TupleData submessage as a key. This field is present if the table in which the delete has happened uses an index as REPLICA IDENTITY.

Byte1('O')

Identifies the following TupleData message as an old tuple. This field is present if the table in which the delete happened has REPLICA IDENTITY set to FULL.

TupleData

TupleData message part representing the contents of the old tuple or primary key, depending on the previous field.

The Delete message may contain either a 'K' message part or an 'O' message part, but never both of them.

Truncate**Byte1('T')**

Identifies the message as a truncate message.

Int64 (TransactionId)

Xid of the transaction (only present for streamed transactions). This field is available since protocol version 2.

Int32

Number of relations

Int8

Option bits for `TRUNCATE`: 1 for `CASCADE`, 2 for `RESTART IDENTITY`

Int32 (Oid)

OID of the relation corresponding to the ID in the relation message. This field is repeated for each relation.

The following messages (Stream Start, Stream Stop, Stream Commit, and Stream Abort) are available since protocol version 2.

Stream Start

Byte1('S')

Identifies the message as a stream start message.

Int64 (TransactionId)

Xid of the transaction.

Int8

A value of 1 indicates this is the first stream segment for this XID, 0 for any other stream segment.

Stream Stop

Byte1('E')

Identifies the message as a stream stop message.

Stream Commit

Byte1('c')

Identifies the message as a stream commit message.

Int64 (TransactionId)

Xid of the transaction.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The LSN of the commit.

Int64 (XLogRecPtr)

The end LSN of the transaction.

Int64 (TimestampTz)

Commit timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Stream Abort

Byte1('A')

Identifies the message as a stream abort message.

Int64 (TransactionId)

Xid of the transaction.

Int64 (TransactionId)

Xid of the subtransaction (will be same as xid of the transaction for top-level transactions).

Int64 (XLogRecPtr)

The LSN of the abort. This field is available since protocol version 4.

Int64 (TimestampTz)

Abort timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01). This field is available since protocol version 4.

The following messages (Begin Prepare, Prepare, Commit Prepared, Rollback Prepared, Stream Prepare) are available since protocol version 3.

Begin Prepare

Byte1('b')

Identifies the message as the beginning of a prepared transaction message.

Int64 (XLogRecPtr)

The LSN of the prepare.

Int64 (XLogRecPtr)

The end LSN of the prepared transaction.

Int64 (TimestampTz)

Prepare timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

String

The user defined GID of the prepared transaction.

Prepare

Byte1('P')

Identifies the message as a prepared transaction message.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The LSN of the prepare.

Int64 (XLogRecPtr)

The end LSN of the prepared transaction.

Int64 (TimestampTz)

Prepare timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

String

The user defined GID of the prepared transaction.

Commit Prepared

Byte1('K')

Identifies the message as the commit of a prepared transaction message.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The LSN of the commit of the prepared transaction.

Int64 (XLogRecPtr)

The end LSN of the commit of the prepared transaction.

Int64 (TimestampTz)

Commit timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

String

The user defined GID of the prepared transaction.

Rollback Prepared

Byte1('r')

Identifies the message as the rollback of a prepared transaction message.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The end LSN of the prepared transaction.

Int64 (XLogRecPtr)

The end LSN of the rollback of the prepared transaction.

Int64 (TimestampTz)

Prepare timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TimestampTz)

Rollback timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

String

The user defined GID of the prepared transaction.

Stream Prepare

Byte1('p')

Identifies the message as a stream prepared transaction message.

Int8(0)

Flags; currently unused.

Int64 (XLogRecPtr)

The LSN of the prepare.

Int64 (XLogRecPtr)

The end LSN of the prepared transaction.

Int64 (TimestampTz)

Prepare timestamp of the transaction. The value is in number of microseconds since PostgreSQL epoch (2000-01-01).

Int64 (TransactionId)

Xid of the transaction.

String

The user defined GID of the prepared transaction.

The following message parts are shared by the above messages.

TupleData

Int16

Number of columns.

Next, one of the following submessages appears for each column (except generated columns):

Byte1('n')

Identifies the data as NULL value.

Or

Byte1('u')

Identifies unchanged TOASTed value (the actual value is not sent).

Or

Byte1('t')

Identifies the data as text formatted value.

Or

Byte1('b')

Identifies the data as binary formatted value.

Int32

Length of the column value.

Byte_{*n*}

The value of the column, either in binary or in text format. (As specified in the preceding format byte). *n* is the above length.

56.10. Summary of Changes since Protocol 2.0

This section provides a quick checklist of changes, for the benefit of developers trying to update existing client libraries to protocol 3.0.

The initial startup packet uses a flexible list-of-strings format instead of a fixed format. Notice that session default values for run-time parameters can now be specified directly in the startup packet. (Actually, you could do that before using the `options` field, but given the limited width of `options` and the lack of any way to quote whitespace in the values, it wasn't a very safe technique.)

All messages now have a length count immediately following the message type byte (except for startup packets, which have no type byte). Also note that `PasswordMessage` now has a type byte.

`ErrorResponse` and `NoticeResponse` ('E' and 'N') messages now contain multiple fields, from which the client code can assemble an error message of the desired level of verbosity. Note that individual fields will typically not end with a newline, whereas the single string sent in the older protocol always did.

The `ReadyForQuery` ('Z') message includes a transaction status indicator.

The distinction between `BinaryRow` and `DataRow` message types is gone; the single `DataRow` message type serves for returning data in all formats. Note that the layout of `DataRow` has changed to make it easier to parse. Also, the representation of binary values has changed: it is no longer directly tied to the server's internal representation.

There is a new “extended query” sub-protocol, which adds the frontend message types `Parse`, `Bind`, `Execute`, `Describe`, `Close`, `Flush`, and `Sync`, and the backend message types `ParseComplete`, `BindComplete`, `PortalSuspended`, `ParameterDescription`, `NoData`, and `CloseComplete`. Existing clients do not have to concern themselves with this sub-protocol, but making use of it might allow improvements in performance or functionality.

`COPY` data is now encapsulated into `CopyData` and `CopyDone` messages. There is a well-defined way to recover from errors during `COPY`. The special “\.” last line is not needed anymore, and is not sent during `COPY OUT`. (It is still recognized as a terminator during `COPY IN`, but its use is deprecated and will eventually be removed.) Binary `COPY` is supported. The `CopyInResponse` and `CopyOutResponse` messages include fields indicating the number of columns and the format of each column.

The layout of `FunctionCall` and `FunctionCallResponse` messages has changed. `FunctionCall` can now support passing `NULL` arguments to functions. It also can handle passing parameters and retrieving results in either text or binary format. There is no longer any reason to consider `FunctionCall` a potential security hole, since it does not offer direct access to internal server data representations.

The backend sends `ParameterStatus` ('S') messages during connection startup for all parameters it considers interesting to the client library. Subsequently, a `ParameterStatus` message is sent whenever the active value changes for any of these parameters.

The `RowDescription` ('T') message carries new table OID and column number fields for each column of the described row. It also shows the format code for each column.

The `CursorResponse` ('P') message is no longer generated by the backend.

The `NotificationResponse` ('A') message has an additional string field, which can carry a “payload” string passed from the `NOTIFY` event sender.

The `EmptyQueryResponse` ('I') message used to include an empty string parameter; this has been removed.

Chapter 57. Writing a Procedural Language Handler

All calls to functions that are written in a language other than the current “version 1” interface for compiled languages (this includes functions in user-defined procedural languages and functions written in SQL) go through a *call handler* function for the specific language. It is the responsibility of the call handler to execute the function in a meaningful way, such as by interpreting the supplied source text. This chapter outlines how a new procedural language's call handler can be written.

The call handler for a procedural language is a “normal” function that must be written in a compiled language such as C, using the version-1 interface, and registered with Postgres Pro as taking no arguments and returning the type `language_handler`. This special pseudo-type identifies the function as a call handler and prevents it from being called directly in SQL commands. For more details on C language calling conventions and dynamic loading, see [Section 39.10](#).

The call handler is called in the same way as any other function: It receives a pointer to a `FunctionCallInfoBaseData` struct containing argument values and information about the called function, and it is expected to return a `Datum` result (and possibly set the `isnull` field of the `FunctionCallInfoBaseData` structure, if it wishes to return an SQL null result). The difference between a call handler and an ordinary callee function is that the `flinfo->fn_oid` field of the `FunctionCallInfoBaseData` structure will contain the OID of the actual function to be called, not of the call handler itself. The call handler must use this field to determine which function to execute. Also, the passed argument list has been set up according to the declaration of the target function, not of the call handler.

It's up to the call handler to fetch the entry of the function from the `pg_proc` system catalog and to analyze the argument and return types of the called function. The `AS` clause from the `CREATE FUNCTION` command for the function will be found in the `prosrc` column of the `pg_proc` row. This is commonly source text in the procedural language, but in theory it could be something else, such as a path name to a file, or anything else that tells the call handler what to do in detail.

Often, the same function is called many times per SQL statement. A call handler can avoid repeated lookups of information about the called function by using the `flinfo->fn_extra` field. This will initially be `NULL`, but can be set by the call handler to point at information about the called function. On subsequent calls, if `flinfo->fn_extra` is already non-`NULL` then it can be used and the information lookup step skipped. The call handler must make sure that `flinfo->fn_extra` is made to point at memory that will live at least until the end of the current query, since an `FmgrInfo` data structure could be kept that long. One way to do this is to allocate the extra data in the memory context specified by `flinfo->fn_mcxt`; such data will normally have the same lifespan as the `FmgrInfo` itself. But the handler could also choose to use a longer-lived memory context so that it can cache function definition information across queries.

When a procedural-language function is invoked as a trigger, no arguments are passed in the usual way, but the `FunctionCallInfoBaseData`'s `context` field points at a `TriggerData` structure, rather than being `NULL` as it is in a plain function call. A language handler should provide mechanisms for procedural-language functions to get at the trigger information.

Although providing a call handler is sufficient to create a minimal procedural language, there are two other functions that can optionally be provided to make the language more convenient to use. These are a *validator* and an *inline handler*. A validator can be provided to allow language-specific checking to be done during [CREATE FUNCTION](#). An inline handler can be provided to allow the language to support anonymous code blocks executed via the [DO](#) command.

If a validator is provided by a procedural language, it must be declared as a function taking a single parameter of type `oid`. The validator's result is ignored, so it is customarily declared to return `void`. The validator will be called at the end of a `CREATE FUNCTION` command that has created or updated a function written in the procedural language. The passed-in OID is the OID of the function's `pg_proc` row. The validator must fetch this row in the usual way, and do whatever checking is appropriate. First, call `CheckFunctionValidatorAccess()` to diagnose explicit calls to the validator that the user could not

achieve through `CREATE FUNCTION`. Typical checks then include verifying that the function's argument and result types are supported by the language, and that the function's body is syntactically correct in the language. If the validator finds the function to be okay, it should just return. If it finds an error, it should report that via the normal `ereport()` error reporting mechanism. Throwing an error will force a transaction rollback and thus prevent the incorrect function definition from being committed.

Validator functions should typically honor the `check_function_bodies` parameter: if it is turned off then any expensive or context-sensitive checking should be skipped. If the language provides for code execution at compilation time, the validator must suppress checks that would induce such execution. In particular, this parameter is turned off by `pg_dump` so that it can load procedural language functions without worrying about side effects or dependencies of the function bodies on other database objects. (Because of this requirement, the call handler should avoid assuming that the validator has fully checked the function. The point of having a validator is not to let the call handler omit checks, but to notify the user immediately if there are obvious errors in a `CREATE FUNCTION` command.) While the choice of exactly what to check is mostly left to the discretion of the validator function, note that the core `CREATE FUNCTION` code only executes `SET` clauses attached to a function when `check_function_bodies` is on. Therefore, checks whose results might be affected by GUC parameters definitely should be skipped when `check_function_bodies` is off, to avoid false failures when restoring a dump.

If an inline handler is provided by a procedural language, it must be declared as a function taking a single parameter of type `internal`. The inline handler's result is ignored, so it is customarily declared to return `void`. The inline handler will be called when a `DO` statement is executed specifying the procedural language. The parameter actually passed is a pointer to an `InlineCodeBlock` struct, which contains information about the `DO` statement's parameters, in particular the text of the anonymous code block to be executed. The inline handler should execute this code and return.

It's recommended that you wrap all these function declarations, as well as the `CREATE LANGUAGE` command itself, into an *extension* so that a simple `CREATE EXTENSION` command is sufficient to install the language. See [Section 39.17](#) for information about writing extensions.

The procedural languages included in the standard distribution are good references when trying to write your own language handler. The [CREATE LANGUAGE](#) reference page has some useful details.

Chapter 58. Writing a Foreign Data Wrapper

All operations on a foreign table are handled through its foreign data wrapper, which consists of a set of functions that the core server calls. The foreign data wrapper is responsible for fetching data from the remote data source and returning it to the Postgres Pro executor. If updating foreign tables is to be supported, the wrapper must handle that, too. This chapter outlines how to write a new foreign data wrapper.

The foreign data wrappers included in the standard distribution are good references when trying to write your own. Look into the `contrib` subdirectory of the source tree.

Note

The SQL standard specifies an interface for writing foreign data wrappers. However, Postgres Pro does not implement that API, because the effort to accommodate it into Postgres Pro would be large, and the standard API hasn't gained wide adoption anyway.

58.1. Foreign Data Wrapper Functions

The FDW author needs to implement a handler function, and optionally a validator function. Both functions must be written in a compiled language such as C, using the version-1 interface. For details on C language calling conventions and dynamic loading, see [Section 39.10](#).

The handler function simply returns a struct of function pointers to callback functions that will be called by the planner, executor, and various maintenance commands. Most of the effort in writing an FDW is in implementing these callback functions. The handler function must be registered with Postgres Pro as taking no arguments and returning the special pseudo-type `fdw_handler`. The callback functions are plain C functions and are not visible or callable at the SQL level. The callback functions are described in [Section 58.2](#).

The validator function is responsible for validating options given in `CREATE` and `ALTER` commands for its foreign data wrapper, as well as foreign servers, user mappings, and foreign tables using the wrapper. The validator function must be registered as taking two arguments, a text array containing the options to be validated, and an OID representing the type of object the options are associated with. The latter corresponds to the OID of the system catalog the object would be stored in, one of:

- `AttributeRelationId`
- `ForeignDataWrapperRelationId`
- `ForeignServerRelationId`
- `ForeignTableRelationId`
- `UserMappingRelationId`

If no validator function is supplied, options are not checked at object creation time or object alteration time.

58.2. Foreign Data Wrapper Callback Routines

The FDW handler function returns a palloc'd `FdwRoutine` struct containing pointers to the callback functions described below. The scan-related functions are required, the rest are optional.

The `FdwRoutine` struct type is declared in [src/include/foreign/fdwapi.h](#), which see for additional details.

58.2.1. FDW Routines for Scanning Foreign Tables

```
void
GetForeignRelSize(PlannerInfo *root,
                  RelOptInfo *baserel,
```

```
Oid foreigntableid);
```

Obtain relation size estimates for a foreign table. This is called at the beginning of planning for a query that scans a foreign table. `root` is the planner's global information about the query; `baserel` is the planner's information about this table; and `foreigntableid` is the `pg_class` OID of the foreign table. (`foreigntableid` could be obtained from the planner data structures, but it's passed explicitly to save effort.)

This function should update `baserel->rows` to be the expected number of rows returned by the table scan, after accounting for the filtering done by the restriction quals. The initial value of `baserel->rows` is just a constant default estimate, which should be replaced if at all possible. The function may also choose to update `baserel->width` if it can compute a better estimate of the average result row width. (The initial value is based on column data types and on column average-width values measured by the last `ANALYZE`.) Also, this function may update `baserel->tuples` if it can compute a better estimate of the foreign table's total row count. (The initial value is from `pg_class.rel tuples` which represents the total row count seen by the last `ANALYZE`; it will be `-1` if no `ANALYZE` has been done on this foreign table.)

See [Section 58.4](#) for additional information.

```
void
GetForeignPaths(PlannerInfo *root,
                RelOptInfo *baserel,
                Oid foreigntableid);
```

Create possible access paths for a scan on a foreign table. This is called during query planning. The parameters are the same as for `GetForeignRelSize`, which has already been called.

This function must generate at least one access path (`ForeignPath` node) for a scan on the foreign table and must call `add_path` to add each such path to `baserel->pathlist`. It's recommended to use `create_foreignscan_path` to build the `ForeignPath` nodes. The function can generate multiple access paths, e.g., a path which has valid `pathkeys` to represent a pre-sorted result. Each access path must contain cost estimates, and can contain any FDW-private information that is needed to identify the specific scan method intended.

See [Section 58.4](#) for additional information.

```
ForeignScan *
GetForeignPlan(PlannerInfo *root,
               RelOptInfo *baserel,
               Oid foreigntableid,
               ForeignPath *best_path,
               List *tlist,
               List *scan_clauses,
               Plan *outer_plan);
```

Create a `ForeignScan` plan node from the selected foreign access path. This is called at the end of query planning. The parameters are as for `GetForeignRelSize`, plus the selected `ForeignPath` (previously produced by `GetForeignPaths`, `GetForeignJoinPaths`, or `GetForeignUpperPaths`), the target list to be emitted by the plan node, the restriction clauses to be enforced by the plan node, and the outer subplan of the `ForeignScan`, which is used for rechecks performed by `RecheckForeignScan`. (If the path is for a join rather than a base relation, `foreigntableid` is `InvalidOid`.)

This function must create and return a `ForeignScan` plan node; it's recommended to use `make_foreignscan` to build the `ForeignScan` node.

See [Section 58.4](#) for additional information.

```
void
BeginForeignScan(ForeignScanState *node,
                 int eflags);
```

Begin executing a foreign scan. This is called during executor startup. It should perform any initialization needed before the scan can start, but not start executing the actual scan (that should be done upon the first call to `IterateForeignScan`). The `ForeignScanState` node has already been created, but its `fdw_state` field is still `NULL`. Information about the table to scan is accessible through the `ForeignScanState` node (in particular, from the underlying `ForeignScan` plan node, which contains any FDW-private information provided by `GetForeignPlan`). `eflags` contains flag bits describing the executor's operating mode for this plan node.

Note that when `(eflags & EXEC_FLAG_EXPLAIN_ONLY)` is true, this function should not perform any externally-visible actions; it should only do the minimum required to make the node state valid for `ExplainForeignScan` and `EndForeignScan`.

```
TupleTableSlot *
IterateForeignScan(ForeignScanState *node);
```

Fetch one row from the foreign source, returning it in a tuple table slot (the node's `ScanTupleSlot` should be used for this purpose). Return `NULL` if no more rows are available. The tuple table slot infrastructure allows either a physical or virtual tuple to be returned; in most cases the latter choice is preferable from a performance standpoint. Note that this is called in a short-lived memory context that will be reset between invocations. Create a memory context in `BeginForeignScan` if you need longer-lived storage, or use the `es_query_cxt` of the node's `EState`.

The rows returned must match the `fdw_scan_tlist` target list if one was supplied, otherwise they must match the row type of the foreign table being scanned. If you choose to optimize away fetching columns that are not needed, you should insert nulls in those column positions, or else generate a `fdw_scan_tlist` list with those columns omitted.

Note that Postgres Pro's executor doesn't care whether the rows returned violate any constraints that were defined on the foreign table — but the planner does care, and may optimize queries incorrectly if there are rows visible in the foreign table that do not satisfy a declared constraint. If a constraint is violated when the user has declared that the constraint should hold true, it may be appropriate to raise an error (just as you would need to do in the case of a data type mismatch).

```
void
ReScanForeignScan(ForeignScanState *node);
```

Restart the scan from the beginning. Note that any parameters the scan depends on may have changed value, so the new scan does not necessarily return exactly the same rows.

```
void
EndForeignScan(ForeignScanState *node);
```

End the scan and release resources. It is normally not important to release `palloc'd` memory, but for example open files and connections to remote servers should be cleaned up.

58.2.2. FDW Routines for Scanning Foreign Joins

If an FDW supports performing foreign joins remotely (rather than by fetching both tables' data and doing the join locally), it should provide this callback function:

```
void
GetForeignJoinPaths(PlannerInfo *root,
                   RelOptInfo *joinrel,
                   RelOptInfo *outerrel,
                   RelOptInfo *innerrel,
                   JoinType jointype,
                   JoinPathExtraData *extra);
```

Create possible access paths for a join of two (or more) foreign tables that all belong to the same foreign server. This optional function is called during query planning. As with `GetForeignPaths`, this function should generate `ForeignPath` path(s) for the supplied `joinrel` (use `create_foreign_join_path` to build

them), and call `add_path` to add these paths to the set of paths considered for the join. But unlike `GetForeignPaths`, it is not necessary that this function succeed in creating at least one path, since paths involving local joining are always possible.

Note that this function will be invoked repeatedly for the same join relation, with different combinations of inner and outer relations; it is the responsibility of the FDW to minimize duplicated work.

Note also that the set of join clauses to apply to the join, which is passed as `extra->restrictlist`, varies depending on the combination of inner and outer relations. A `ForeignPath` path generated for the `joinrel` must contain the set of join clauses it uses, which will be used by the planner to convert the `ForeignPath` path into a plan, if it is selected by the planner as the best path for the `joinrel`.

If a `ForeignPath` path is chosen for the join, it will represent the entire join process; paths generated for the component tables and subsidiary joins will not be used. Subsequent processing of the join path proceeds much as it does for a path scanning a single foreign table. One difference is that the `scanrelid` of the resulting `ForeignScan` plan node should be set to zero, since there is no single relation that it represents; instead, the `fs_relids` field of the `ForeignScan` node represents the set of relations that were joined. (The latter field is set up automatically by the core planner code, and need not be filled by the FDW.) Another difference is that, because the column list for a remote join cannot be found from the system catalogs, the FDW must fill `fdw_scan_tlist` with an appropriate list of `TargetEntry` nodes, representing the set of columns it will supply at run time in the tuples it returns.

Note

Beginning with Postgres Pro 16, `fs_relids` includes the rangetable indexes of outer joins, if any were involved in this join. The new field `fs_base_relids` includes only base relation indexes, and thus mimics `fs_relids`'s old semantics.

See [Section 58.4](#) for additional information.

58.2.3. FDW Routines for Planning Post-Scan/Join Processing

If an FDW supports performing remote post-scan/join processing, such as remote aggregation, it should provide this callback function:

```
void
GetForeignUpperPaths(PlannerInfo *root,
                    UpperRelationKind stage,
                    RelOptInfo *input_rel,
                    RelOptInfo *output_rel,
                    void *extra);
```

Create possible access paths for *upper relation* processing, which is the planner's term for all post-scan/join query processing, such as aggregation, window functions, sorting, and table updates. This optional function is called during query planning. Currently, it is called only if all base relation(s) involved in the query belong to the same FDW. This function should generate `ForeignPath` path(s) for any post-scan/join processing that the FDW knows how to perform remotely (use `create_foreign_upper_path` to build them), and call `add_path` to add these paths to the indicated upper relation. As with `GetForeignJoinPaths`, it is not necessary that this function succeed in creating any paths, since paths involving local processing are always possible.

The `stage` parameter identifies which post-scan/join step is currently being considered. `output_rel` is the upper relation that should receive paths representing computation of this step, and `input_rel` is the relation representing the input to this step. The `extra` parameter provides additional details, currently, it is set only for `UPPERREL_PARTIAL_GROUP_AGG` or `UPPERREL_GROUP_AGG`, in which case it points to a `GroupPathExtraData` structure; or for `UPPERREL_FINAL`, in which case it points to a `FinalPathExtraData` structure. (Note that `ForeignPath` paths added to `output_rel` would typically not have any direct dependency on paths of the `input_rel`, since their processing is expected to be done externally.

However, examining paths previously generated for the previous processing step can be useful to avoid redundant planning work.)

See [Section 58.4](#) for additional information.

58.2.4. FDW Routines for Updating Foreign Tables

If an FDW supports writable foreign tables, it should provide some or all of the following callback functions depending on the needs and capabilities of the FDW:

```
void
AddForeignUpdateTargets(PlannerInfo *root,
                        Index rtindex,
                        RangeTblEntry *target_rte,
                        Relation target_relation);
```

UPDATE and DELETE operations are performed against rows previously fetched by the table-scanning functions. The FDW may need extra information, such as a row ID or the values of primary-key columns, to ensure that it can identify the exact row to update or delete. To support that, this function can add extra hidden, or “junk”, target columns to the list of columns that are to be retrieved from the foreign table during an UPDATE or DELETE.

To do that, construct a Var representing an extra value you need, and pass it to `add_row_identity_var`, along with a name for the junk column. (You can do this more than once if several columns are needed.) You must choose a distinct junk column name for each different Var you need, except that Vars that are identical except for the `varno` field can and should share a column name. The core system uses the junk column names `tableoid` for a table's `tableoid` column, `ctid` or `ctidN` for `ctid`, `wholerow` for a whole-row Var marked with `vartype = RECORD`, and `wholerowN` for a whole-row Var with `vartype` equal to the table's declared row type. Re-use these names when you can (the planner will combine duplicate requests for identical junk columns). If you need another kind of junk column besides these, it might be wise to choose a name prefixed with your extension name, to avoid conflicts against other FDWs.

If the `AddForeignUpdateTargets` pointer is set to NULL, no extra target expressions are added. (This will make it impossible to implement DELETE operations, though UPDATE may still be feasible if the FDW relies on an unchanging primary key to identify rows.)

```
List *
PlanForeignModify(PlannerInfo *root,
                  ModifyTable *plan,
                  Index resultRelation,
                  int subplan_index);
```

Perform any additional planning actions needed for an insert, update, or delete on a foreign table. This function generates the FDW-private information that will be attached to the `ModifyTable` plan node that performs the update action. This private information must have the form of a List, and will be delivered to `BeginForeignModify` during the execution stage.

`root` is the planner's global information about the query. `plan` is the `ModifyTable` plan node, which is complete except for the `fdwPrivLists` field. `resultRelation` identifies the target foreign table by its range table index. `subplan_index` identifies which target of the `ModifyTable` plan node this is, counting from zero; use this if you want to index into per-target-relation substructures of the `plan` node.

See [Section 58.4](#) for additional information.

If the `PlanForeignModify` pointer is set to NULL, no additional plan-time actions are taken, and the `fdw_private` list delivered to `BeginForeignModify` will be NIL.

```
void
BeginForeignModify(ModifyTableState *mtstate,
                  ResultRelInfo *rinfo,
                  List *fdw_private,
```



```
int subplan_index,  
int eflags);
```

Begin executing a foreign table modification operation. This routine is called during executor startup. It should perform any initialization needed prior to the actual table modifications. Subsequently, `ExecForeignInsert`/`ExecForeignBatchInsert`, `ExecForeignUpdate` or `ExecForeignDelete` will be called for tuple(s) to be inserted, updated, or deleted.

`mtstate` is the overall state of the `ModifyTable` plan node being executed; global data about the plan and execution state is available via this structure. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. (The `ri_FdwState` field of `ResultRelInfo` is available for the FDW to store any private state it needs for this operation.) `fdw_private` contains the private data generated by `PlanForeignModify`, if any. `subplan_index` identifies which target of the `ModifyTable` plan node this is. `eflags` contains flag bits describing the executor's operating mode for this plan node.

Note that when `(eflags & EXEC_FLAG_EXPLAIN_ONLY)` is true, this function should not perform any externally-visible actions; it should only do the minimum required to make the node state valid for `ExplainForeignModify` and `EndForeignModify`.

If the `BeginForeignModify` pointer is set to `NULL`, no action is taken during executor startup.

```
TupleTableSlot *  
ExecForeignInsert(EState *estate,  
                  ResultRelInfo *rinfo,  
                  TupleTableSlot *slot,  
                  TupleTableSlot *planSlot);
```

Insert one tuple into the foreign table. `estate` is global execution state for the query. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. `slot` contains the tuple to be inserted; it will match the row-type definition of the foreign table. `planSlot` contains the tuple that was generated by the `ModifyTable` plan node's subplan; it differs from `slot` in possibly containing additional "junk" columns. (The `planSlot` is typically of little interest for `INSERT` cases, but is provided for completeness.)

The return value is either a slot containing the data that was actually inserted (this might differ from the data supplied, for example as a result of trigger actions), or `NULL` if no row was actually inserted (again, typically as a result of triggers). The passed-in `slot` can be re-used for this purpose.

The data in the returned slot is used only if the `INSERT` statement has a `RETURNING` clause or involves a view `WITH CHECK OPTION`; or if the foreign table has an `AFTER ROW` trigger. Triggers require all columns, but the FDW could choose to optimize away returning some or all columns depending on the contents of the `RETURNING` clause or `WITH CHECK OPTION` constraints. Regardless, some slot must be returned to indicate success, or the query's reported row count will be wrong.

If the `ExecForeignInsert` pointer is set to `NULL`, attempts to insert into the foreign table will fail with an error message.

Note that this function is also called when inserting routed tuples into a foreign-table partition or executing `COPY FROM` on a foreign table, in which case it is called in a different way than it is in the `INSERT` case. See the callback functions described below that allow the FDW to support that.

```
TupleTableSlot **  
ExecForeignBatchInsert(EState *estate,  
                      ResultRelInfo *rinfo,  
                      TupleTableSlot **slots,  
                      TupleTableSlot **planSlots,  
                      int *numSlots);
```

Insert multiple tuples in bulk into the foreign table. The parameters are the same for `ExecForeignInsert` except `slots` and `planSlots` contain multiple tuples and `*numSlots` specifies the number of tuples in those arrays.

The return value is an array of slots containing the data that was actually inserted (this might differ from the data supplied, for example as a result of trigger actions.) The passed-in `slots` can be re-used for this purpose. The number of successfully inserted tuples is returned in `*numSlots`.

The data in the returned slot is used only if the `INSERT` statement involves a view `WITH CHECK OPTION`; or if the foreign table has an `AFTER ROW` trigger. Triggers require all columns, but the FDW could choose to optimize away returning some or all columns depending on the contents of the `WITH CHECK OPTION` constraints.

If the `ExecForeignBatchInsert` or `GetForeignModifyBatchSize` pointer is set to `NULL`, attempts to insert into the foreign table will use `ExecForeignInsert`. This function is not used if the `INSERT` has the `RETURNING` clause.

Note that this function is also called when inserting routed tuples into a foreign-table partition or executing `COPY FROM` on a foreign table, in which case it is called in a different way than it is in the `INSERT` case. See the callback functions described below that allow the FDW to support that.

```
int
GetForeignModifyBatchSize(ResultRelInfo *rinfo);
```

Report the maximum number of tuples that a single `ExecForeignBatchInsert` call can handle for the specified foreign table. The executor passes at most the given number of tuples to `ExecForeignBatchInsert`. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. The FDW is expected to provide a foreign server and/or foreign table option for the user to set this value, or some hard-coded value.

If the `ExecForeignBatchInsert` or `GetForeignModifyBatchSize` pointer is set to `NULL`, attempts to insert into the foreign table will use `ExecForeignInsert`.

```
TupleTableSlot *
ExecForeignUpdate(EState *estate,
                  ResultRelInfo *rinfo,
                  TupleTableSlot *slot,
                  TupleTableSlot *planSlot);
```

Update one tuple in the foreign table. `estate` is global execution state for the query. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. `slot` contains the new data for the tuple; it will match the row-type definition of the foreign table. `planSlot` contains the tuple that was generated by the `ModifyTable` plan node's subplan. Unlike `slot`, this tuple contains only the new values for columns changed by the query, so do not rely on attribute numbers of the foreign table to index into `planSlot`. Also, `planSlot` typically contains additional "junk" columns. In particular, any junk columns that were requested by `AddForeignUpdateTargets` will be available from this slot.

The return value is either a slot containing the row as it was actually updated (this might differ from the data supplied, for example as a result of trigger actions), or `NULL` if no row was actually updated (again, typically as a result of triggers). The passed-in `slot` can be re-used for this purpose.

The data in the returned slot is used only if the `UPDATE` statement has a `RETURNING` clause or involves a view `WITH CHECK OPTION`; or if the foreign table has an `AFTER ROW` trigger. Triggers require all columns, but the FDW could choose to optimize away returning some or all columns depending on the contents of the `RETURNING` clause or `WITH CHECK OPTION` constraints. Regardless, some slot must be returned to indicate success, or the query's reported row count will be wrong.

If the `ExecForeignUpdate` pointer is set to `NULL`, attempts to update the foreign table will fail with an error message.

```
TupleTableSlot *
ExecForeignDelete(EState *estate,
                  ResultRelInfo *rinfo,
                  TupleTableSlot *slot,
```

```
TupleTableSlot *planSlot);
```

Delete one tuple from the foreign table. `estate` is global execution state for the query. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. `slot` contains nothing useful upon call, but can be used to hold the returned tuple. `planSlot` contains the tuple that was generated by the `ModifyTable` plan node's subplan; in particular, it will carry any junk columns that were requested by `AddForeignUpdateTargets`. The junk column(s) must be used to identify the tuple to be deleted.

The return value is either a slot containing the row that was deleted, or `NULL` if no row was deleted (typically as a result of triggers). The passed-in `slot` can be used to hold the tuple to be returned.

The data in the returned slot is used only if the `DELETE` query has a `RETURNING` clause or the foreign table has an `AFTER ROW` trigger. Triggers require all columns, but the FDW could choose to optimize away returning some or all columns depending on the contents of the `RETURNING` clause. Regardless, some slot must be returned to indicate success, or the query's reported row count will be wrong.

If the `ExecForeignDelete` pointer is set to `NULL`, attempts to delete from the foreign table will fail with an error message.

```
void
EndForeignModify (EState *estate,
                  ResultRelInfo *rinfo);
```

End the table update and release resources. It is normally not important to release `palloc'd` memory, but for example open files and connections to remote servers should be cleaned up.

If the `EndForeignModify` pointer is set to `NULL`, no action is taken during executor shutdown.

Tuples inserted into a partitioned table by `INSERT` or `COPY FROM` are routed to partitions. If an FDW supports routable foreign-table partitions, it should also provide the following callback functions. These functions are also called when `COPY FROM` is executed on a foreign table.

```
void
BeginForeignInsert (ModifyTableState *mtstate,
                   ResultRelInfo *rinfo);
```

Begin executing an insert operation on a foreign table. This routine is called right before the first tuple is inserted into the foreign table in both cases when it is the partition chosen for tuple routing and the target specified in a `COPY FROM` command. It should perform any initialization needed prior to the actual insertion. Subsequently, `ExecForeignInsert` or `ExecForeignBatchInsert` will be called for tuple(s) to be inserted into the foreign table.

`mtstate` is the overall state of the `ModifyTable` plan node being executed; global data about the plan and execution state is available via this structure. `rinfo` is the `ResultRelInfo` struct describing the target foreign table. (The `ri_FdwState` field of `ResultRelInfo` is available for the FDW to store any private state it needs for this operation.)

When this is called by a `COPY FROM` command, the plan-related global data in `mtstate` is not provided and the `planSlot` parameter of `ExecForeignInsert` subsequently called for each inserted tuple is `NULL`, whether the foreign table is the partition chosen for tuple routing or the target specified in the command.

If the `BeginForeignInsert` pointer is set to `NULL`, no action is taken for the initialization.

Note that if the FDW does not support routable foreign-table partitions and/or executing `COPY FROM` on foreign tables, this function or `ExecForeignInsert/ExecForeignBatchInsert` subsequently called must throw error as needed.

```
void
EndForeignInsert (EState *estate,
                  ResultRelInfo *rinfo);
```

End the insert operation and release resources. It is normally not important to release `palloc'd` memory, but for example open files and connections to remote servers should be cleaned up.

If the `EndForeignInsert` pointer is set to `NULL`, no action is taken for the termination.

```
int
IsForeignRelUpdatable(Relation rel);
```

Report which update operations the specified foreign table supports. The return value should be a bit mask of rule event numbers indicating which operations are supported by the foreign table, using the `CmdType` enumeration; that is, $(1 \ll \text{CMD_UPDATE}) = 4$ for `UPDATE`, $(1 \ll \text{CMD_INSERT}) = 8$ for `INSERT`, and $(1 \ll \text{CMD_DELETE}) = 16$ for `DELETE`.

If the `IsForeignRelUpdatable` pointer is set to `NULL`, foreign tables are assumed to be insertable, updatable, or deletable if the FDW provides `ExecForeignInsert`, `ExecForeignUpdate`, or `ExecForeignDelete` respectively. This function is only needed if the FDW supports some tables that are updatable and some that are not. (Even then, it's permissible to throw an error in the execution routine instead of checking in this function. However, this function is used to determine updatability for display in the `information_schema` views.)

Some inserts, updates, and deletes to foreign tables can be optimized by implementing an alternative set of interfaces. The ordinary interfaces for inserts, updates, and deletes fetch rows from the remote server and then modify those rows one at a time. In some cases, this row-by-row approach is necessary, but it can be inefficient. If it is possible for the foreign server to determine which rows should be modified without actually retrieving them, and if there are no local structures which would affect the operation (row-level local triggers, stored generated columns, or `WITH CHECK OPTION` constraints from parent views), then it is possible to arrange things so that the entire operation is performed on the remote server. The interfaces described below make this possible.

```
bool
PlanDirectModify(PlannerInfo *root,
                 ModifyTable *plan,
                 Index resultRelation,
                 int subplan_index);
```

Decide whether it is safe to execute a direct modification on the remote server. If so, return `true` after performing planning actions needed for that. Otherwise, return `false`. This optional function is called during query planning. If this function succeeds, `BeginDirectModify`, `IterateDirectModify` and `EndDirectModify` will be called at the execution stage, instead. Otherwise, the table modification will be executed using the table-updating functions described above. The parameters are the same as for `PlanForeignModify`.

To execute the direct modification on the remote server, this function must rewrite the target subplan with a `ForeignScan` plan node that executes the direct modification on the remote server. The `operation` and `resultRelation` fields of the `ForeignScan` must be set appropriately. `operation` must be set to the `CmdType` enumeration corresponding to the statement kind (that is, `CMD_UPDATE` for `UPDATE`, `CMD_INSERT` for `INSERT`, and `CMD_DELETE` for `DELETE`), and the `resultRelation` argument must be copied to the `resultRelation` field.

See [Section 58.4](#) for additional information.

If the `PlanDirectModify` pointer is set to `NULL`, no attempts to execute a direct modification on the remote server are taken.

```
void
BeginDirectModify(ForeignScanState *node,
                 int eflags);
```

Prepare to execute a direct modification on the remote server. This is called during executor startup. It should perform any initialization needed prior to the direct modification (that should be done upon

the first call to `IterateDirectModify`). The `ForeignScanState` node has already been created, but its `fdw_state` field is still `NULL`. Information about the table to modify is accessible through the `ForeignScanState` node (in particular, from the underlying `ForeignScan` plan node, which contains any FDW-private information provided by `PlanDirectModify`). `eflags` contains flag bits describing the executor's operating mode for this plan node.

Note that when `(eflags & EXEC_FLAG_EXPLAIN_ONLY)` is true, this function should not perform any externally-visible actions; it should only do the minimum required to make the node state valid for `ExplainDirectModify` and `EndDirectModify`.

If the `BeginDirectModify` pointer is set to `NULL`, no attempts to execute a direct modification on the remote server are taken.

```
TupleTableSlot *
IterateDirectModify(ForeignScanState *node);
```

When the `INSERT`, `UPDATE` or `DELETE` query doesn't have a `RETURNING` clause, just return `NULL` after a direct modification on the remote server. When the query has the clause, fetch one result containing the data needed for the `RETURNING` calculation, returning it in a tuple table slot (the node's `ScanTupleSlot` should be used for this purpose). The data that was actually inserted, updated or deleted must be stored in `node->resultRelInfo->ri_projectReturning->pi_exprContext->ecxt_scantuple`. Return `NULL` if no more rows are available. Note that this is called in a short-lived memory context that will be reset between invocations. Create a memory context in `BeginDirectModify` if you need longer-lived storage, or use the `es_query_cxt` of the node's `EState`.

The rows returned must match the `fdw_scan_tlist` target list if one was supplied, otherwise they must match the row type of the foreign table being updated. If you choose to optimize away fetching columns that are not needed for the `RETURNING` calculation, you should insert nulls in those column positions, or else generate a `fdw_scan_tlist` list with those columns omitted.

Whether the query has the clause or not, the query's reported row count must be incremented by the FDW itself. When the query doesn't have the clause, the FDW must also increment the row count for the `ForeignScanState` node in the `EXPLAIN ANALYZE` case.

If the `IterateDirectModify` pointer is set to `NULL`, no attempts to execute a direct modification on the remote server are taken.

```
void
EndDirectModify(ForeignScanState *node);
```

Clean up following a direct modification on the remote server. It is normally not important to release `palloc'd` memory, but for example open files and connections to the remote server should be cleaned up.

If the `EndDirectModify` pointer is set to `NULL`, no attempts to execute a direct modification on the remote server are taken.

58.2.5. FDW Routines for `TRUNCATE`

```
void
ExecForeignTruncate(List *rels,
                    DropBehavior behavior,
                    bool restart_seqs);
```

Truncate foreign tables. This function is called when `TRUNCATE` is executed on a foreign table. `rels` is a list of `Relation` data structures of foreign tables to truncate.

`behavior` is either `DROP_RESTRICT` or `DROP_CASCADE` indicating that the `RESTRICT` or `CASCADE` option was requested in the original `TRUNCATE` command, respectively.

If `restart_seqs` is true, the original `TRUNCATE` command requested the `RESTART IDENTITY` behavior, otherwise the `CONTINUE IDENTITY` behavior was requested.

Note that the `ONLY` options specified in the original `TRUNCATE` command are not passed to `ExecForeignTruncate`. This behavior is similar to the callback functions of `SELECT`, `UPDATE` and `DELETE` on a foreign table.

`ExecForeignTruncate` is invoked once per foreign server for which foreign tables are to be truncated. This means that all foreign tables included in `rels` must belong to the same server.

If the `ExecForeignTruncate` pointer is set to `NULL`, attempts to truncate foreign tables will fail with an error message.

58.2.6. FDW Routines for Row Locking

If an FDW wishes to support *late row locking* (as described in [Section 58.5](#)), it must provide the following callback functions:

```
RowMarkType
GetForeignRowMarkType(RangeTblEntry *rte,
                      LockClauseStrength strength);
```

Report which row-marking option to use for a foreign table. `rte` is the `RangeTblEntry` node for the table and `strength` describes the lock strength requested by the relevant `FOR UPDATE/SHARE` clause, if any. The result must be a member of the `RowMarkType` enum type.

This function is called during query planning for each foreign table that appears in an `UPDATE`, `DELETE`, or `SELECT FOR UPDATE/SHARE` query and is not the target of `UPDATE` or `DELETE`.

If the `GetForeignRowMarkType` pointer is set to `NULL`, the `ROW_MARK_COPY` option is always used. (This implies that `RefetchForeignRow` will never be called, so it need not be provided either.)

See [Section 58.5](#) for more information.

```
void
RefetchForeignRow(EState *estate,
                  ExecRowMark *erm,
                  Datum rowid,
                  TupleTableSlot *slot,
                  bool *updated);
```

Re-fetch one tuple slot from the foreign table, after locking it if required. `estate` is global execution state for the query. `erm` is the `ExecRowMark` struct describing the target foreign table and the row lock type (if any) to acquire. `rowid` identifies the tuple to be fetched. `slot` contains nothing useful upon call, but can be used to hold the returned tuple. `updated` is an output parameter.

This function should store the tuple into the provided slot, or clear it if the row lock couldn't be obtained. The row lock type to acquire is defined by `erm->markType`, which is the value previously returned by `GetForeignRowMarkType`. (`ROW_MARK_REFERENCE` means to just re-fetch the tuple without acquiring any lock, and `ROW_MARK_COPY` will never be seen by this routine.)

In addition, `*updated` should be set to `true` if what was fetched was an updated version of the tuple rather than the same version previously obtained. (If the FDW cannot be sure about this, always returning `true` is recommended.)

Note that by default, failure to acquire a row lock should result in raising an error; returning with an empty slot is only appropriate if the `SKIP LOCKED` option is specified by `erm->waitPolicy`.

The `rowid` is the `ctid` value previously read for the row to be re-fetched. Although the `rowid` value is passed as a `Datum`, it can currently only be a `tid`. The function API is chosen in hopes that it may be possible to allow other data types for row IDs in future.

If the `RefetchForeignRow` pointer is set to `NULL`, attempts to re-fetch rows will fail with an error message.

See [Section 58.5](#) for more information.

```
bool
RecheckForeignScan(ForeignScanState *node,
                   TupleTableSlot *slot);
```

Recheck that a previously-returned tuple still matches the relevant scan and join qualifiers, and possibly provide a modified version of the tuple. For foreign data wrappers which do not perform join pushdown, it will typically be more convenient to set this to `NULL` and instead set `fdw_recheck_qual`s appropriately. When outer joins are pushed down, however, it isn't sufficient to reapply the checks relevant to all the base tables to the result tuple, even if all needed attributes are present, because failure to match some qualifier might result in some attributes going to `NULL`, rather than in no tuple being returned. `RecheckForeignScan` can recheck qualifiers and return true if they are still satisfied and false otherwise, but it can also store a replacement tuple into the supplied slot.

To implement join pushdown, a foreign data wrapper will typically construct an alternative local join plan which is used only for rechecks; this will become the outer subplan of the `ForeignScan`. When a recheck is required, this subplan can be executed and the resulting tuple can be stored in the slot. This plan need not be efficient since no base table will return more than one row; for example, it may implement all joins as nested loops. The function `GetExistingLocalJoinPath` may be used to search existing paths for a suitable local join path, which can be used as the alternative local join plan. `GetExistingLocalJoinPath` searches for an unparameterized path in the path list of the specified join relation. (If it does not find such a path, it returns `NULL`, in which case a foreign data wrapper may build the local path by itself or may choose not to create access paths for that join.)

58.2.7. FDW Routines for EXPLAIN

```
void
ExplainForeignScan(ForeignScanState *node,
                   ExplainState *es);
```

Print additional `EXPLAIN` output for a foreign table scan. This function can call `ExplainPropertyText` and related functions to add fields to the `EXPLAIN` output. The flag fields in `es` can be used to determine what to print, and the state of the `ForeignScanState` node can be inspected to provide run-time statistics in the `EXPLAIN ANALYZE` case.

If the `ExplainForeignScan` pointer is set to `NULL`, no additional information is printed during `EXPLAIN`.

```
void
ExplainForeignModify(ModifyTableState *mtstate,
                    ResultRelInfo *rinfo,
                    List *fdw_private,
                    int subplan_index,
                    struct ExplainState *es);
```

Print additional `EXPLAIN` output for a foreign table update. This function can call `ExplainPropertyText` and related functions to add fields to the `EXPLAIN` output. The flag fields in `es` can be used to determine what to print, and the state of the `ModifyTableState` node can be inspected to provide run-time statistics in the `EXPLAIN ANALYZE` case. The first four arguments are the same as for `BeginForeignModify`.

If the `ExplainForeignModify` pointer is set to `NULL`, no additional information is printed during `EXPLAIN`.

```
void
ExplainDirectModify(ForeignScanState *node,
                   ExplainState *es);
```

Print additional `EXPLAIN` output for a direct modification on the remote server. This function can call `ExplainPropertyText` and related functions to add fields to the `EXPLAIN` output. The flag fields in `es` can be used to determine what to print, and the state of the `ForeignScanState` node can be inspected to provide run-time statistics in the `EXPLAIN ANALYZE` case.

If the `ExplainDirectModify` pointer is set to `NULL`, no additional information is printed during `EXPLAIN`.

58.2.8. FDW Routines for ANALYZE

```
bool
AnalyzeForeignTable(Relation relation,
                    AcquireSampleRowsFunc *func,
                    BlockNumber *totalpages);
```

This function is called when **ANALYZE** is executed on a foreign table. If the FDW can collect statistics for this foreign table, it should return `true`, and provide a pointer to a function that will collect sample rows from the table in `func`, plus the estimated size of the table in pages in `totalpages`. Otherwise, return `false`.

If the FDW does not support collecting statistics for any tables, the `AnalyzeForeignTable` pointer can be set to `NULL`.

If provided, the sample collection function must have the signature

```
int
AcquireSampleRowsFunc(Relation relation,
                      int elevel,
                      HeapTuple *rows,
                      int targrows,
                      double *totalrows,
                      double *totaldeadrows);
```

A random sample of up to `targrows` rows should be collected from the table and stored into the caller-provided `rows` array. The actual number of rows collected must be returned. In addition, store estimates of the total numbers of live and dead rows in the table into the output parameters `totalrows` and `totaldeadrows`. (Set `totaldeadrows` to zero if the FDW does not have any concept of dead rows.)

58.2.9. FDW Routines for Parallel Execution

A `ForeignScan` node can, optionally, support parallel execution. A parallel `ForeignScan` will be executed in multiple processes and must return each row exactly once across all cooperating processes. To do this, processes can coordinate through fixed-size chunks of dynamic shared memory. This shared memory is not guaranteed to be mapped at the same address in every process, so it must not contain pointers. The following functions are all optional, but most are required if parallel execution is to be supported.

```
bool
IsForeignScanParallelSafe(PlannerInfo *root, RelOptInfo *rel,
                          RangeTblEntry *rte);
```

Test whether a scan can be performed within a parallel worker. This function will only be called when the planner believes that a parallel plan might be possible, and should return `true` if it is safe for that scan to run within a parallel worker. This will generally not be the case if the remote data source has transaction semantics, unless the worker's connection to the data can somehow be made to share the same transaction context as the leader.

If this function is not defined, it is assumed that the scan must take place within the parallel leader. Note that returning `true` does not mean that the scan itself can be done in parallel, only that the scan can be performed within a parallel worker. Therefore, it can be useful to define this method even when parallel execution is not supported.

```
Size
EstimateDSMForeignScan(ForeignScanState *node, ParallelContext *pcxt);
```

Estimate the amount of dynamic shared memory that will be required for parallel operation. This may be higher than the amount that will actually be used, but it must not be lower. The return value is in bytes. This function is optional, and can be omitted if not needed; but if it is omitted, the next three functions must be omitted as well, because no shared memory will be allocated for the FDW's use.


```
void
InitializeDSMForeignScan(ForeignScanState *node, ParallelContext *pcxt,
                        void *coordinate);
```

Initialize the dynamic shared memory that will be required for parallel operation. `coordinate` points to a shared memory area of size equal to the return value of `EstimateDSMForeignScan`. This function is optional, and can be omitted if not needed.

```
void
ReInitializeDSMForeignScan(ForeignScanState *node, ParallelContext *pcxt,
                          void *coordinate);
```

Re-initialize the dynamic shared memory required for parallel operation when the foreign-scan plan node is about to be re-scanned. This function is optional, and can be omitted if not needed. Recommended practice is that this function reset only shared state, while the `ReScanForeignScan` function resets only local state. Currently, this function will be called before `ReScanForeignScan`, but it's best not to rely on that ordering.

```
void
InitializeWorkerForeignScan(ForeignScanState *node, shm_toc *toc,
                          void *coordinate);
```

Initialize a parallel worker's local state based on the shared state set up by the leader during `InitializeDSMForeignScan`. This function is optional, and can be omitted if not needed.

```
void
ShutdownForeignScan(ForeignScanState *node);
```

Release resources when it is anticipated the node will not be executed to completion. This is not called in all cases; sometimes, `EndForeignScan` may be called without this function having been called first. Since the DSM segment used by parallel query is destroyed just after this callback is invoked, foreign data wrappers that wish to take some action before the DSM segment goes away should implement this method.

58.2.10. FDW Routines for Asynchronous Execution

A `ForeignScan` node can, optionally, support asynchronous execution. The following functions are all optional, but are all required if asynchronous execution is to be supported.

```
bool
IsForeignPathAsyncCapable(ForeignPath *path);
```

Test whether a given `ForeignPath` `path` can scan the underlying foreign relation asynchronously. This function will only be called at the end of query planning when the given path is a direct child of an `AppendPath` path and when the planner believes that asynchronous execution improves performance, and should return true if the given path is able to scan the foreign relation asynchronously.

If this function is not defined, it is assumed that the given path scans the foreign relation using `IterateForeignScan`. (This implies that the callback functions described below will never be called, so they need not be provided either.)

```
void
ForeignAsyncRequest(AsyncRequest *areq);
```

Produce one tuple asynchronously from the `ForeignScan` node. `areq` is the `AsyncRequest` struct describing the `ForeignScan` node and the parent `Append` node that requested the tuple from it. This function should store the tuple into the slot specified by `areq->result`, and set `areq->request_complete` to true; or if it needs to wait on an event external to the core server such as network I/O, and cannot produce any tuple immediately, set the flag to false, and set `areq->callback_pending` to true for the `ForeignScan` node to get a callback from the callback functions described below. If no more tuples are available, set the slot to NULL or an empty slot, and the `areq->request_complete` flag to true. It's recommended to use `ExecAsyncRequestDone` or `ExecAsyncRequestPending` to set the output parameters in the `areq`.

```
void
ForeignAsyncConfigureWait (AsyncRequest *areq);
```

Configure a file descriptor event for which the `ForeignScan` node wishes to wait. This function will only be called when the `ForeignScan` node has the `areq->callback_pending` flag set, and should add the event to the `as_eventset` of the parent `Append` node described by the `areq`. When the file descriptor event occurs, `ForeignAsyncNotify` will be called.

```
void
ForeignAsyncNotify (AsyncRequest *areq);
```

Process a relevant event that has occurred, then produce one tuple asynchronously from the `ForeignScan` node. This function should set the output parameters in the `areq` in the same way as `ForeignAsyncRequest`.

58.2.11. FDW Routines for Reparameterization of Paths

```
List *
ReparameterizeForeignPathByChild (PlannerInfo *root, List *fdw_private,
                                  RelOptInfo *child_rel);
```

This function is called while converting a path parameterized by the top-most parent of the given child relation `child_rel` to be parameterized by the child relation. The function is used to reparameterize any paths or translate any expression nodes saved in the given `fdw_private` member of a `ForeignPath`. The callback may use `reparameterize_path_by_child`, `adjust_appendrel_attrs` or `adjust_appendrel_attrs_multilevel` as required.

58.3. Foreign Data Wrapper Helper Functions

Several helper functions are exported from the core server so that authors of foreign data wrappers can get easy access to attributes of FDW-related objects, such as FDW options. To use any of these functions, you need to include the header file `foreign/foreign.h` in your source file. That header also defines the struct types that are returned by these functions.

```
ForeignDataWrapper *
GetForeignDataWrapperExtended (Oid fdwid, bits16 flags);
```

This function returns a `ForeignDataWrapper` object for the foreign-data wrapper with the given OID. A `ForeignDataWrapper` object contains properties of the FDW (see `foreign/foreign.h` for details). `flags` is a bitwise-or'd bit mask indicating an extra set of options. It can take the value `FDW_MISSING_OK`, in which case a `NULL` result is returned to the caller instead of an error for an undefined object.

```
ForeignDataWrapper *
GetForeignDataWrapper (Oid fdwid);
```

This function returns a `ForeignDataWrapper` object for the foreign-data wrapper with the given OID. A `ForeignDataWrapper` object contains properties of the FDW (see `foreign/foreign.h` for details).

```
ForeignServer *
GetForeignServerExtended (Oid serverid, bits16 flags);
```

This function returns a `ForeignServer` object for the foreign server with the given OID. A `ForeignServer` object contains properties of the server (see `foreign/foreign.h` for details). `flags` is a bitwise-or'd bit mask indicating an extra set of options. It can take the value `FSV_MISSING_OK`, in which case a `NULL` result is returned to the caller instead of an error for an undefined object.

```
ForeignServer *
GetForeignServer (Oid serverid);
```

This function returns a `ForeignServer` object for the foreign server with the given OID. A `ForeignServer` object contains properties of the server (see `foreign/foreign.h` for details).

```
UserMapping *
```

```
GetUserMapping(Oid userid, Oid serverid);
```

This function returns a `UserMapping` object for the user mapping of the given role on the given server. (If there is no mapping for the specific user, it will return the mapping for `PUBLIC`, or throw error if there is none.) A `UserMapping` object contains properties of the user mapping (see `foreign/foreign.h` for details).

```
ForeignTable *  
GetForeignTable(Oid relid);
```

This function returns a `ForeignTable` object for the foreign table with the given OID. A `ForeignTable` object contains properties of the foreign table (see `foreign/foreign.h` for details).

```
List *  
GetForeignColumnOptions(Oid relid, AttrNumber attnum);
```

This function returns the per-column FDW options for the column with the given foreign table OID and attribute number, in the form of a list of `DefElem`. NIL is returned if the column has no options.

Some object types have name-based lookup functions in addition to the OID-based ones:

```
ForeignDataWrapper *  
GetForeignDataWrapperByName(const char *name, bool missing_ok);
```

This function returns a `ForeignDataWrapper` object for the foreign-data wrapper with the given name. If the wrapper is not found, return NULL if `missing_ok` is true, otherwise raise an error.

```
ForeignServer *  
GetForeignServerByName(const char *name, bool missing_ok);
```

This function returns a `ForeignServer` object for the foreign server with the given name. If the server is not found, return NULL if `missing_ok` is true, otherwise raise an error.

58.4. Foreign Data Wrapper Query Planning

The FDW callback functions `GetForeignRelSize`, `GetForeignPaths`, `GetForeignPlan`, `PlanForeignModify`, `GetForeignJoinPaths`, `GetForeignUpperPaths`, and `PlanDirectModify` must fit into the workings of the Postgres Pro planner. Here are some notes about what they must do.

The information in `root` and `baserel` can be used to reduce the amount of information that has to be fetched from the foreign table (and therefore reduce the cost). `baserel->baserestrictinfo` is particularly interesting, as it contains restriction quals (`WHERE` clauses) that should be used to filter the rows to be fetched. (The FDW itself is not required to enforce these quals, as the core executor can check them instead.) `baserel->reldata->exprs` can be used to determine which columns need to be fetched; but note that it only lists columns that have to be emitted by the `ForeignScan` plan node, not columns that are used in qual evaluation but not output by the query.

Various private fields are available for the FDW planning functions to keep information in. Generally, whatever you store in FDW private fields should be `palloc'd`, so that it will be reclaimed at the end of planning.

`baserel->fdw_private` is a void pointer that is available for FDW planning functions to store information relevant to the particular foreign table. The core planner does not touch it except to initialize it to NULL when the `RelOptInfo` node is created. It is useful for passing information forward from `GetForeignRelSize` to `GetForeignPaths` and/or `GetForeignPaths` to `GetForeignPlan`, thereby avoiding recalculation.

`GetForeignPaths` can identify the meaning of different access paths by storing private information in the `fdw_private` field of `ForeignPath` nodes. `fdw_private` is declared as a `List` pointer, but could actually contain anything since the core planner does not touch it. However, best practice is to use a representation that's dumpable by `nodeToString`, for use with debugging support available in the backend.

`GetForeignPlan` can examine the `fdw_private` field of the selected `ForeignPath` node, and can generate `fdw_exprs` and `fdw_private` lists to be placed in the `ForeignScan` plan node, where they will be available at execution time. Both of these lists must be represented in a form that `copyObject` knows how to copy. The `fdw_private` list has no other restrictions and is not interpreted by the core backend in any way. The `fdw_exprs` list, if not `NIL`, is expected to contain expression trees that are intended to be executed at run time. These trees will undergo post-processing by the planner to make them fully executable.

In `GetForeignPlan`, generally the passed-in target list can be copied into the plan node as-is. The passed `scan_clauses` list contains the same clauses as `baserel->baserestrictinfo`, but may be re-ordered for better execution efficiency. In simple cases the FDW can just strip `RestrictInfo` nodes from the `scan_clauses` list (using `extract_actual_clauses`) and put all the clauses into the plan node's qual list, which means that all the clauses will be checked by the executor at run time. More complex FDWs may be able to check some of the clauses internally, in which case those clauses can be removed from the plan node's qual list so that the executor doesn't waste time rechecking them.

As an example, the FDW might identify some restriction clauses of the form `foreign_variable = sub_expression`, which it determines can be executed on the remote server given the locally-evaluated value of the `sub_expression`. The actual identification of such a clause should happen during `GetForeignPaths`, since it would affect the cost estimate for the path. The path's `fdw_private` field would probably include a pointer to the identified clause's `RestrictInfo` node. Then `GetForeignPlan` would remove that clause from `scan_clauses`, but add the `sub_expression` to `fdw_exprs` to ensure that it gets massaged into executable form. It would probably also put control information into the plan node's `fdw_private` field to tell the execution functions what to do at run time. The query transmitted to the remote server would involve something like `WHERE foreign_variable = $1`, with the parameter value obtained at run time from evaluation of the `fdw_exprs` expression tree.

Any clauses removed from the plan node's qual list must instead be added to `fdw_recheckquals` or rechecked by `RecheckForeignScan` in order to ensure correct behavior at the `READ COMMITTED` isolation level. When a concurrent update occurs for some other table involved in the query, the executor may need to verify that all of the original quals are still satisfied for the tuple, possibly against a different set of parameter values. Using `fdw_recheckquals` is typically easier than implementing checks inside `RecheckForeignScan`, but this method will be insufficient when outer joins have been pushed down, since the join tuples in that case might have some fields go to `NULL` without rejecting the tuple entirely.

Another `ForeignScan` field that can be filled by FDWs is `fdw_scan_tlist`, which describes the tuples returned by the FDW for this plan node. For simple foreign table scans this can be set to `NIL`, implying that the returned tuples have the row type declared for the foreign table. A non-`NIL` value must be a target list (list of `TargetEntry`s) containing Vars and/or expressions representing the returned columns. This might be used, for example, to show that the FDW has omitted some columns that it noticed won't be needed for the query. Also, if the FDW can compute expressions used by the query more cheaply than can be done locally, it could add those expressions to `fdw_scan_tlist`. Note that join plans (created from paths made by `GetForeignJoinPaths`) must always supply `fdw_scan_tlist` to describe the set of columns they will return.

The FDW should always construct at least one path that depends only on the table's restriction clauses. In join queries, it might also choose to construct path(s) that depend on join clauses, for example `foreign_variable = local_variable`. Such clauses will not be found in `baserel->baserestrictinfo` but must be sought in the relation's join lists. A path using such a clause is called a "parameterized path". It must identify the other relations used in the selected join clause(s) with a suitable value of `param_info`; use `get_baserel_parampathinfo` to compute that value. In `GetForeignPlan`, the `local_variable` portion of the join clause would be added to `fdw_exprs`, and then at run time the case works the same as for an ordinary restriction clause.

If an FDW supports remote joins, `GetForeignJoinPaths` should produce `ForeignPaths` for potential remote joins in much the same way as `GetForeignPaths` works for base tables. Information about the intended join can be passed forward to `GetForeignPlan` in the same ways described above. However, `baserestrictinfo` is not relevant for join relations; instead, the relevant join clauses for a particular join are passed to `GetForeignJoinPaths` as a separate parameter (`extra->restrictlist`).

An FDW might additionally support direct execution of some plan actions that are above the level of scans and joins, such as grouping or aggregation. To offer such options, the FDW should generate paths and insert them into the appropriate *upper relation*. For example, a path representing remote aggregation should be inserted into the `UPPERREL_GROUP_AGG` relation, using `add_path`. This path will be compared on a cost basis with local aggregation performed by reading a simple scan path for the foreign relation (note that such a path must also be supplied, else there will be an error at plan time). If the remote-aggregation path wins, which it usually would, it will be converted into a plan in the usual way, by calling `GetForeignPlan`. The recommended place to generate such paths is in the `GetForeignUpperPaths` callback function, which is called for each upper relation (i.e., each post-scan/join processing step), if all the base relations of the query come from the same FDW.

`PlanForeignModify` and the other callbacks described in [Section 58.2.4](#) are designed around the assumption that the foreign relation will be scanned in the usual way and then individual row updates will be driven by a local `ModifyTable` plan node. This approach is necessary for the general case where an update requires reading local tables as well as foreign tables. However, if the operation could be executed entirely by the foreign server, the FDW could generate a path representing that and insert it into the `UPPERREL_FINAL` upper relation, where it would compete against the `ModifyTable` approach. This approach could also be used to implement remote `SELECT FOR UPDATE`, rather than using the row locking callbacks described in [Section 58.2.6](#). Keep in mind that a path inserted into `UPPERREL_FINAL` is responsible for implementing *all* behavior of the query.

When planning an `UPDATE` or `DELETE`, `PlanForeignModify` and `PlanDirectModify` can look up the `RelOptInfo` struct for the foreign table and make use of the `baserel->fdw_private` data previously created by the scan-planning functions. However, in `INSERT` the target table is not scanned so there is no `RelOptInfo` for it. The `List` returned by `PlanForeignModify` has the same restrictions as the `fdw_private` list of a `ForeignScan` plan node, that is it must contain only structures that `copyObject` knows how to copy.

`INSERT` with an `ON CONFLICT` clause does not support specifying the conflict target, as unique constraints or exclusion constraints on remote tables are not locally known. This in turn implies that `ON CONFLICT DO UPDATE` is not supported, since the specification is mandatory there.

58.5. Row Locking in Foreign Data Wrappers

If an FDW's underlying storage mechanism has a concept of locking individual rows to prevent concurrent updates of those rows, it is usually worthwhile for the FDW to perform row-level locking with as close an approximation as practical to the semantics used in ordinary Postgres Pro tables. There are multiple considerations involved in this.

One key decision to be made is whether to perform *early locking* or *late locking*. In early locking, a row is locked when it is first retrieved from the underlying store, while in late locking, the row is locked only when it is known that it needs to be locked. (The difference arises because some rows may be discarded by locally-checked restriction or join conditions.) Early locking is much simpler and avoids extra round trips to a remote store, but it can cause locking of rows that need not have been locked, resulting in reduced concurrency or even unexpected deadlocks. Also, late locking is only possible if the row to be locked can be uniquely re-identified later. Preferably the row identifier should identify a specific version of the row, as Postgres Pro TIDs do.

By default, Postgres Pro ignores locking considerations when interfacing to FDWs, but an FDW can perform early locking without any explicit support from the core code. The API functions described in [Section 58.2.6](#), which were added in Postgres Pro 9.5, allow an FDW to use late locking if it wishes.

An additional consideration is that in `READ COMMITTED` isolation mode, Postgres Pro may need to re-check restriction and join conditions against an updated version of some target tuple. Rechecking join conditions requires re-obtaining copies of the non-target rows that were previously joined to the target tuple. When working with standard Postgres Pro tables, this is done by including the TIDs of the non-target tables in the column list projected through the join, and then re-fetching non-target rows when required. This approach keeps the join data set compact, but it requires inexpensive re-fetch capability, as well as a TID that can uniquely identify the row version to be re-fetched. By default, therefore, the approach used with foreign tables is to include a copy of the entire row fetched from a foreign table in

the column list projected through the join. This puts no special demands on the FDW but can result in reduced performance of merge and hash joins. An FDW that is capable of meeting the re-fetch requirements can choose to do it the first way.

For an `UPDATE` or `DELETE` on a foreign table, it is recommended that the `ForeignScan` operation on the target table perform early locking on the rows that it fetches, perhaps via the equivalent of `SELECT FOR UPDATE`. An FDW can detect whether a table is an `UPDATE/DELETE` target at plan time by comparing its `relid` to `root->parse->resultRelation`, or at execution time by using `ExecRelationIsTargetRelation()`. An alternative possibility is to perform late locking within the `ExecForeignUpdate` or `ExecForeignDelete` callback, but no special support is provided for this.

For foreign tables that are specified to be locked by a `SELECT FOR UPDATE/SHARE` command, the `ForeignScan` operation can again perform early locking by fetching tuples with the equivalent of `SELECT FOR UPDATE/SHARE`. To perform late locking instead, provide the callback functions defined in [Section 58.2.6](#). In `GetForeignRowMarkType`, select rowmark option `ROW_MARK_EXCLUSIVE`, `ROW_MARK_NOKEYEXCLUSIVE`, `ROW_MARK_SHARE`, or `ROW_MARK_KEYSHARE` depending on the requested lock strength. (The core code will act the same regardless of which of these four options you choose.) Elsewhere, you can detect whether a foreign table was specified to be locked by this type of command by using `get_plan_rowmark` at plan time, or `ExecFindRowMark` at execution time; you must check not only whether a non-null rowmark struct is returned, but that its `strength` field is not `LCS_NONE`.

Lastly, for foreign tables that are used in an `UPDATE`, `DELETE` or `SELECT FOR UPDATE/SHARE` command but are not specified to be row-locked, you can override the default choice to copy entire rows by having `GetForeignRowMarkType` select option `ROW_MARK_REFERENCE` when it sees lock strength `LCS_NONE`. This will cause `RefetchForeignRow` to be called with that value for `markType`; it should then re-fetch the row without acquiring any new lock. (If you have a `GetForeignRowMarkType` function but don't wish to re-fetch unlocked rows, select option `ROW_MARK_COPY` for `LCS_NONE`.)

Chapter 59. Writing a Table Sampling Method

Postgres Pro's implementation of the `TABLESAMPLE` clause supports custom table sampling methods, in addition to the `BERNOULLI` and `SYSTEM` methods that are required by the SQL standard. The sampling method determines which rows of the table will be selected when the `TABLESAMPLE` clause is used.

At the SQL level, a table sampling method is represented by a single SQL function, typically implemented in C, having the signature

```
method_name(internal) RETURNS tsm_handler
```

The name of the function is the same method name appearing in the `TABLESAMPLE` clause. The `internal` argument is a dummy (always having value zero) that simply serves to prevent this function from being called directly from an SQL command. The result of the function must be a palloc'd struct of type `TsmRoutine`, which contains pointers to support functions for the sampling method. These support functions are plain C functions and are not visible or callable at the SQL level. The support functions are described in [Section 59.1](#).

In addition to function pointers, the `TsmRoutine` struct must provide these additional fields:

```
List *parameterTypes
```

This is an OID list containing the data type OIDs of the parameter(s) that will be accepted by the `TABLESAMPLE` clause when this sampling method is used. For example, for the built-in methods, this list contains a single item with value `FLOAT4OID`, which represents the sampling percentage. Custom sampling methods can have more or different parameters.

```
bool repeatable_across_queries
```

If `true`, the sampling method can deliver identical samples across successive queries, if the same parameters and `REPEATABLE` seed value are supplied each time and the table contents have not changed. When this is `false`, the `REPEATABLE` clause is not accepted for use with the sampling method.

```
bool repeatable_across_scans
```

If `true`, the sampling method can deliver identical samples across successive scans in the same query (assuming unchanging parameters, seed value, and snapshot). When this is `false`, the planner will not select plans that would require scanning the sampled table more than once, since that might result in inconsistent query output.

The table sampling methods included in the standard distribution are good references when trying to write your own. Look into the `contrib` subdirectory for add-on methods.

59.1. Sampling Method Support Functions

The TSM handler function returns a palloc'd `TsmRoutine` struct containing pointers to the support functions described below. Most of the functions are required, but some are optional, and those pointers can be `NULL`.

```
void
SampleScanGetSampleSize (PlannerInfo *root,
                        RelOptInfo *baserel,
                        List *paramexprs,
                        BlockNumber *pages,
                        double *tuples);
```

This function is called during planning. It must estimate the number of relation pages that will be read during a sample scan, and the number of tuples that will be selected by the scan. (For example, these

might be determined by estimating the sampling fraction, and then multiplying the `basere1->pages` and `basere1->tuples` numbers by that, being sure to round the results to integral values.) The `paramexprs` list holds the expression(s) that are parameters to the `TABLESAMPLE` clause. It is recommended to use `estimate_expression_value()` to try to reduce these expressions to constants, if their values are needed for estimation purposes; but the function must provide size estimates even if they cannot be reduced, and it should not fail even if the values appear invalid (remember that they're only estimates of what the run-time values will be). The `pages` and `tuples` parameters are outputs.

```
void
InitSampleScan (SampleScanState *node,
                int eflags);
```

Initialize for execution of a `SampleScan` plan node. This is called during executor startup. It should perform any initialization needed before processing can start. The `SampleScanState` node has already been created, but its `tsm_state` field is `NULL`. The `InitSampleScan` function can `malloc` whatever internal state data is needed by the sampling method, and store a pointer to it in `node->tsm_state`. Information about the table to scan is accessible through other fields of the `SampleScanState` node (but note that the `node->ss.ss_currentScanDesc` scan descriptor is not set up yet). `eflags` contains flag bits describing the executor's operating mode for this plan node.

When `(eflags & EXEC_FLAG_EXPLAIN_ONLY)` is true, the scan will not actually be performed, so this function should only do the minimum required to make the node state valid for `EXPLAIN` and `EndSampleScan`.

This function can be omitted (set the pointer to `NULL`), in which case `BeginSampleScan` must perform all initialization needed by the sampling method.

```
void
BeginSampleScan (SampleScanState *node,
                 Datum *params,
                 int nparams,
                 uint32 seed);
```

Begin execution of a sampling scan. This is called just before the first attempt to fetch a tuple, and may be called again if the scan needs to be restarted. Information about the table to scan is accessible through fields of the `SampleScanState` node (but note that the `node->ss.ss_currentScanDesc` scan descriptor is not set up yet). The `params` array, of length `nparams`, contains the values of the parameters supplied in the `TABLESAMPLE` clause. These will have the number and types specified in the sampling method's `parameterTypes` list, and have been checked to not be null. `seed` contains a seed to use for any random numbers generated within the sampling method; it is either a hash derived from the `REPEATABLE` value if one was given, or the result of `random()` if not.

This function may adjust the fields `node->use_bulkread` and `node->use_pagemode`. If `node->use_bulkread` is true, which it is by default, the scan will use a buffer access strategy that encourages recycling buffers after use. It might be reasonable to set this to false if the scan will visit only a small fraction of the table's pages. If `node->use_pagemode` is true, which it is by default, the scan will perform visibility checking in a single pass for all tuples on each visited page. It might be reasonable to set this to false if the scan will select only a small fraction of the tuples on each visited page. That will result in fewer tuple visibility checks being performed, though each one will be more expensive because it will require more locking.

If the sampling method is marked `repeatable_across_scans`, it must be able to select the same set of tuples during a rescan as it did originally, that is a fresh call of `BeginSampleScan` must lead to selecting the same tuples as before (if the `TABLESAMPLE` parameters and seed don't change).

```
BlockNumber
NextSampleBlock (SampleScanState *node, BlockNumber nblocks);
```

Returns the block number of the next page to be scanned, or `InvalidBlockNumber` if no pages remain to be scanned.

This function can be omitted (set the pointer to NULL), in which case the core code will perform a sequential scan of the entire relation. Such a scan can use synchronized scanning, so that the sampling method cannot assume that the relation pages are visited in the same order on each scan.

```
OffsetNumber
NextSampleTuple (SampleScanState *node,
                 BlockNumber blockno,
                 OffsetNumber maxoffset);
```

Returns the offset number of the next tuple to be sampled on the specified page, or `InvalidOffsetNumber` if no tuples remain to be sampled. `maxoffset` is the largest offset number in use on the page.

Note

`NextSampleTuple` is not explicitly told which of the offset numbers in the range `1 .. maxoffset` actually contain valid tuples. This is not normally a problem since the core code ignores requests to sample missing or invisible tuples; that should not result in any bias in the sample. However, if necessary, the function can use `node->donetuples` to examine how many of the tuples it returned were valid and visible.

Note

`NextSampleTuple` must *not* assume that `blockno` is the same page number returned by the most recent `NextSampleBlock` call. It was returned by some previous `NextSampleBlock` call, but the core code is allowed to call `NextSampleBlock` in advance of actually scanning pages, so as to support prefetching. It is OK to assume that once sampling of a given page begins, successive `NextSampleTuple` calls all refer to the same page until `InvalidOffsetNumber` is returned.

```
void
EndSampleScan (SampleScanState *node);
```

End the scan and release resources. It is normally not important to release palloc'd memory, but any externally-visible resources should be cleaned up. This function can be omitted (set the pointer to NULL) in the common case where no such resources exist.

Chapter 60. Writing a Custom Scan Provider

Postgres Pro supports a set of experimental facilities which are intended to allow extension modules to add new scan types to the system. Unlike a [foreign data wrapper](#), which is only responsible for knowing how to scan its own foreign tables, a custom scan provider can provide an alternative method of scanning any relation in the system. Typically, the motivation for writing a custom scan provider will be to allow the use of some optimization not supported by the core system, such as caching or some form of hardware acceleration. This chapter outlines how to write a new custom scan provider.

Implementing a new type of custom scan is a three-step process. First, during planning, it is necessary to generate access paths representing a scan using the proposed strategy. Second, if one of those access paths is selected by the planner as the optimal strategy for scanning a particular relation, the access path must be converted to a plan. Finally, it must be possible to execute the plan and generate the same results that would have been generated for any other access path targeting the same relation.

60.1. Creating Custom Scan Paths

A custom scan provider will typically add paths for a base relation by setting the following hook, which is called after the core code has generated all the access paths it can for the relation (except for Gather and Gather Merge paths, which are made after this call so that they can use partial paths added by the hook):

```
typedef void (*set_rel_pathlist_hook_type) (PlannerInfo *root,
                                           RelOptInfo *rel,
                                           Index rti,
                                           RangeTblEntry *rte);

extern PGDLLIMPORT set_rel_pathlist_hook_type set_rel_pathlist_hook;
```

Although this hook function can be used to examine, modify, or remove paths generated by the core system, a custom scan provider will typically confine itself to generating `CustomPath` objects and adding them to `rel` using `add_path`, or `add_partial_path` if they are partial paths. The custom scan provider is responsible for initializing the `CustomPath` object, which is declared like this:

```
typedef struct CustomPath
{
    Path      path;
    uint32    flags;
    List      *custom_paths;
    List      *custom_restrictinfo;
    List      *custom_private;
    const CustomPathMethods *methods;
} CustomPath;
```

`path` must be initialized as for any other path, including the row-count estimate, start and total cost, and sort ordering provided by this path. `flags` is a bit mask, which specifies whether the scan provider can support certain optional capabilities. `flags` should include `CUSTOMPATH_SUPPORT_BACKWARD_SCAN` if the custom path can support a backward scan, `CUSTOMPATH_SUPPORT_MARK_RESTORE` if it can support mark and restore, and `CUSTOMPATH_SUPPORT_PROJECTION` if it can perform projections. (If `CUSTOMPATH_SUPPORT_PROJECTION` is not set, the scan node will only be asked to produce Vars of the scanned relation; while if that flag is set, the scan node must be able to evaluate scalar expressions over these Vars.) An optional `custom_paths` is a list of `Path` nodes used by this custom-path node; these will be transformed into `Plan` nodes by planner. As described below, custom paths can be created for join relations as well. In such a case, `custom_restrictinfo` should be used to store the set of join clauses to apply to the join the custom path replaces. Otherwise it should be `NIL`. `custom_private` can be used to store the custom path's private data. Private data should be stored in a form that can be handled by `nodeToString`, so that debugging routines that attempt to print the custom path will work as designed. `methods` must point to a (usually statically allocated) object implementing the required custom path methods, which are further detailed below.

A custom scan provider can also provide join paths. Just as for base relations, such a path must produce the same output as would normally be produced by the join it replaces. To do this, the join provider should set the following hook, and then within the hook function, create `CustomPath` path(s) for the join relation.

```
typedef void (*set_join_pathlist_hook_type) (PlannerInfo *root,
                                             RelOptInfo *joinrel,
                                             RelOptInfo *outerrel,
                                             RelOptInfo *innerrel,
                                             JoinType jointype,
                                             JoinPathExtraData *extra);

extern PGDLLIMPORT set_join_pathlist_hook_type set_join_pathlist_hook;
```

This hook will be invoked repeatedly for the same join relation, with different combinations of inner and outer relations; it is the responsibility of the hook to minimize duplicated work.

Note also that the set of join clauses to apply to the join, which is passed as `extra->restrictlist`, varies depending on the combination of inner and outer relations. A `CustomPath` path generated for the `joinrel` must contain the set of join clauses it uses, which will be used by the planner to convert the `CustomPath` path into a plan, if it is selected by the planner as the best path for the `joinrel`.

60.1.1. Custom Scan Path Callbacks

```
Plan *(*PlanCustomPath) (PlannerInfo *root,
                         RelOptInfo *rel,
                         CustomPath *best_path,
                         List *tlist,
                         List *clauses,
                         List *custom_plans);
```

Convert a custom path to a finished plan. The return value will generally be a `CustomScan` object, which the callback must allocate and initialize. See [Section 60.2](#) for more details.

```
List *(*ReparameterizeCustomPathByChild) (PlannerInfo *root,
                                           List *custom_private,
                                           RelOptInfo *child_rel);
```

This callback is called while converting a path parameterized by the top-most parent of the given child relation `child_rel` to be parameterized by the child relation. The callback is used to reparameterize any paths or translate any expression nodes saved in the given `custom_private` member of a `CustomPath`. The callback may use `reparameterize_path_by_child`, `adjust_appendrel_attrs` or `adjust_appendrel_attrs_multilevel` as required.

60.2. Creating Custom Scan Plans

A custom scan is represented in a finished plan tree using the following structure:

```
typedef struct CustomScan
{
    Scan        scan;
    uint32      flags;
    List        *custom_plans;
    List        *custom_exprs;
    List        *custom_private;
    List        *custom_scan_tlist;
    Bitmapset   *custom_relids;
    const CustomScanMethods *methods;
} CustomScan;
```

`scan` must be initialized as for any other scan, including estimated costs, target lists, qualifications, and so on. `flags` is a bit mask with the same meaning as in `CustomPath`. `custom_plans` can be used to store child `Plan` nodes. `custom_exprs` should be used to store expression trees that will need to be fixed up

by `setrefs.c` and `subselect.c`, while `custom_private` should be used to store other private data that is only used by the custom scan provider itself. `custom_scan_tlist` can be `NIL` when scanning a base relation, indicating that the custom scan returns scan tuples that match the base relation's row type. Otherwise it is a target list describing the actual scan tuples. `custom_scan_tlist` must be provided for joins, and could be provided for scans if the custom scan provider can compute some non-Var expressions. `custom_relids` is set by the core code to the set of relations (range table indexes) that this scan node handles; except when this scan is replacing a join, it will have only one member. `methods` must point to a (usually statically allocated) object implementing the required custom scan methods, which are further detailed below.

When a `CustomScan` scans a single relation, `scan.scanrelid` must be the range table index of the table to be scanned. When it replaces a join, `scan.scanrelid` should be zero.

Plan trees must be able to be duplicated using `copyObject`, so all the data stored within the “custom” fields must consist of nodes that that function can handle. Furthermore, custom scan providers cannot substitute a larger structure that embeds a `CustomScan` for the structure itself, as would be possible for a `CustomPath` or `CustomScanState`.

60.2.1. Custom Scan Plan Callbacks

```
Node *(*CreateCustomScanState) (CustomScan *cscan);
```

Allocate a `CustomScanState` for this `CustomScan`. The actual allocation will often be larger than required for an ordinary `CustomScanState`, because many providers will wish to embed that as the first field of a larger structure. The value returned must have the node tag and `methods` set appropriately, but other fields should be left as zeroes at this stage; after `ExecInitCustomScan` performs basic initialization, the `BeginCustomScan` callback will be invoked to give the custom scan provider a chance to do whatever else is needed.

60.3. Executing Custom Scans

When a `CustomScan` is executed, its execution state is represented by a `CustomScanState`, which is declared as follows:

```
typedef struct CustomScanState
{
    ScanState ss;
    uint32    flags;
    const CustomExecMethods *methods;
} CustomScanState;
```

`ss` is initialized as for any other scan state, except that if the scan is for a join rather than a base relation, `ss.ss_currentRelation` is left `NULL`. `flags` is a bit mask with the same meaning as in `CustomPath` and `CustomScan`. `methods` must point to a (usually statically allocated) object implementing the required custom scan state methods, which are further detailed below. Typically, a `CustomScanState`, which need not support `copyObject`, will actually be a larger structure embedding the above as its first member.

60.3.1. Custom Scan Execution Callbacks

```
void (*BeginCustomScan) (CustomScanState *node,
                        EState *estate,
                        int eflags);
```

Complete initialization of the supplied `CustomScanState`. Standard fields have been initialized by `ExecInitCustomScan`, but any private fields should be initialized here.

```
TupleTableSlot *(*ExecCustomScan) (CustomScanState *node);
```

Fetch the next scan tuple. If any tuples remain, it should fill `ps_ResultTupleSlot` with the next tuple in the current scan direction, and then return the tuple slot. If not, `NULL` or an empty slot should be returned.

```
void (*EndCustomScan) (CustomScanState *node);
```

Clean up any private data associated with the `CustomScanState`. This method is required, but it does not need to do anything if there is no associated data or it will be cleaned up automatically.

```
void (*ReScanCustomScan) (CustomScanState *node);
```

Rewind the current scan to the beginning and prepare to rescan the relation.

```
void (*MarkPosCustomScan) (CustomScanState *node);
```

Save the current scan position so that it can subsequently be restored by the `RestrPosCustomScan` callback. This callback is optional, and need only be supplied if the `CUSTOMPATH_SUPPORT_MARK_RESTORE` flag is set.

```
void (*RestrPosCustomScan) (CustomScanState *node);
```

Restore the previous scan position as saved by the `MarkPosCustomScan` callback. This callback is optional, and need only be supplied if the `CUSTOMPATH_SUPPORT_MARK_RESTORE` flag is set.

```
Size (*EstimateDSMCustomScan) (CustomScanState *node,  
                               ParallelContext *pcxt);
```

Estimate the amount of dynamic shared memory that will be required for parallel operation. This may be higher than the amount that will actually be used, but it must not be lower. The return value is in bytes. This callback is optional, and need only be supplied if this custom scan provider supports parallel execution.

```
void (*InitializeDSMCustomScan) (CustomScanState *node,  
                                ParallelContext *pcxt,  
                                void *coordinate);
```

Initialize the dynamic shared memory that will be required for parallel operation. `coordinate` points to a shared memory area of size equal to the return value of `EstimateDSMCustomScan`. This callback is optional, and need only be supplied if this custom scan provider supports parallel execution.

```
void (*ReInitializeDSMCustomScan) (CustomScanState *node,  
                                  ParallelContext *pcxt,  
                                  void *coordinate);
```

Re-initialize the dynamic shared memory required for parallel operation when the custom-scan plan node is about to be re-scanned. This callback is optional, and need only be supplied if this custom scan provider supports parallel execution. Recommended practice is that this callback reset only shared state, while the `ReScanCustomScan` callback resets only local state. Currently, this callback will be called before `ReScanCustomScan`, but it's best not to rely on that ordering.

```
void (*InitializeWorkerCustomScan) (CustomScanState *node,  
                                    shm_toc *toc,  
                                    void *coordinate);
```

Initialize a parallel worker's local state based on the shared state set up by the leader during `InitializeDSMCustomScan`. This callback is optional, and need only be supplied if this custom scan provider supports parallel execution.

```
void (*ShutdownCustomScan) (CustomScanState *node);
```

Release resources when it is anticipated the node will not be executed to completion. This is not called in all cases; sometimes, `EndCustomScan` may be called without this function having been called first. Since the DSM segment used by parallel query is destroyed just after this callback is invoked, custom scan providers that wish to take some action before the DSM segment goes away should implement this method.

```
void (*ExplainCustomScan) (CustomScanState *node,  
                           List *ancestors,
```

```
ExplainState *es);
```

Output additional information for `EXPLAIN` of a custom-scan plan node. This callback is optional. Common data stored in the `ScanState`, such as the target list and scan relation, will be shown even without this callback, but the callback allows the display of additional, private state.

Chapter 61. Genetic Query Optimizer

Author

Written by Martin Utesch (<utesch@aut.tu-freiberg.de>) for the Institute of Automatic Control at the University of Mining and Technology in Freiberg, Germany.

61.1. Query Handling as a Complex Optimization Problem

Among all relational operators the most difficult one to process and optimize is the *join*. The number of possible query plans grows exponentially with the number of joins in the query. Further optimization effort is caused by the support of a variety of *join methods* (e.g., nested loop, hash join, merge join in Postgres Pro) to process individual joins and a diversity of *indexes* (e.g., B-tree, hash, GiST and GIN in Postgres Pro) as access paths for relations.

The normal Postgres Pro query optimizer performs a *near-exhaustive search* over the space of alternative strategies. This algorithm, first introduced in IBM's System R database, produces a near-optimal join order, but can take an enormous amount of time and memory space when the number of joins in the query grows large. This makes the ordinary Postgres Pro query optimizer inappropriate for queries that join a large number of tables.

The Institute of Automatic Control at the University of Mining and Technology, in Freiberg, Germany, encountered some problems when it wanted to use PostgreSQL as the backend for a decision support knowledge based system for the maintenance of an electrical power grid. The DBMS needed to handle large join queries for the inference machine of the knowledge based system. The number of joins in these queries made using the normal query optimizer infeasible.

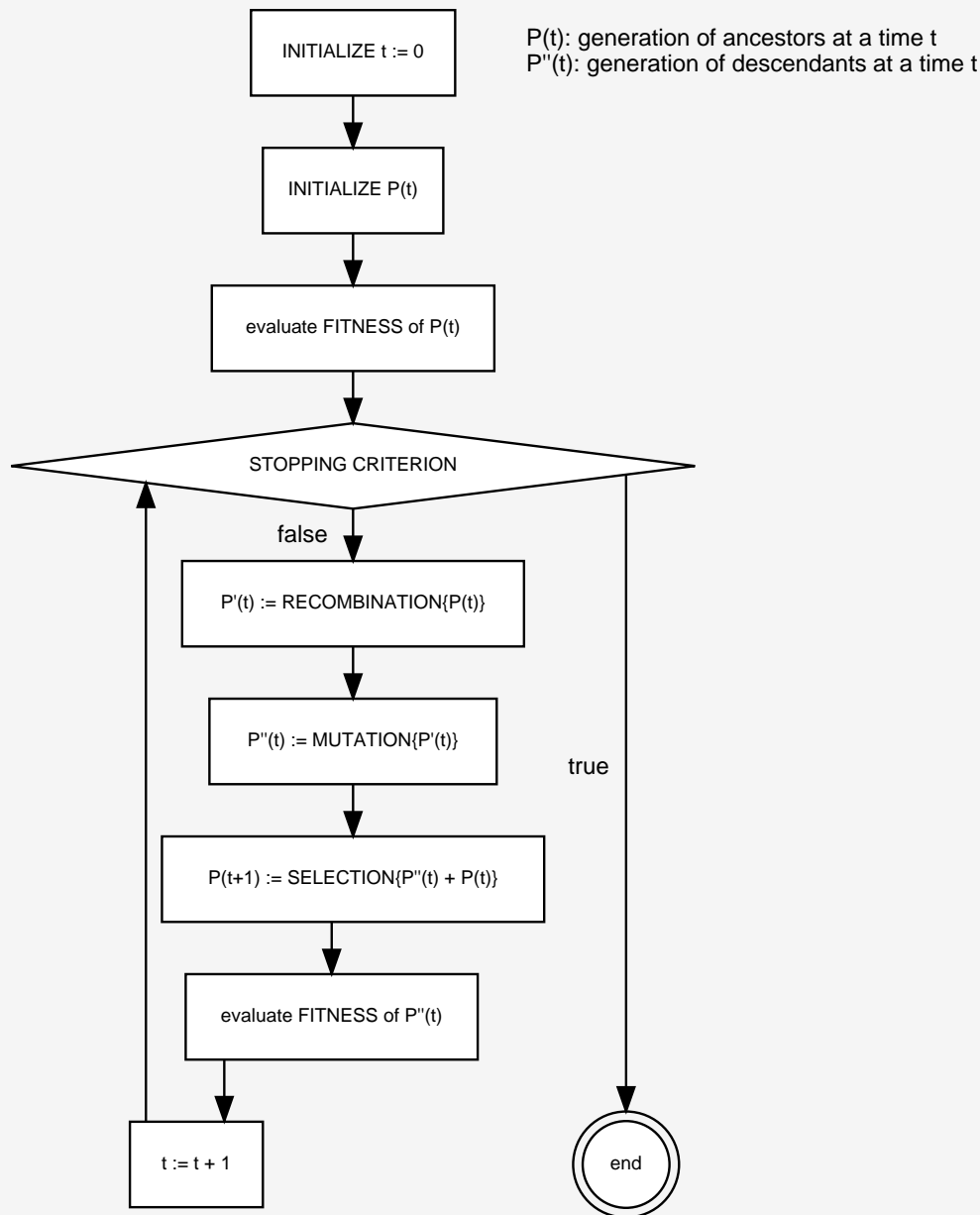
In the following we describe the implementation of a *genetic algorithm* to solve the join ordering problem in a manner that is efficient for queries involving large numbers of joins.

61.2. Genetic Algorithms

The genetic algorithm (GA) is a heuristic optimization method which operates through randomized search. The set of possible solutions for the optimization problem is considered as a *population of individuals*. The degree of adaptation of an individual to its environment is specified by its *fitness*.

The coordinates of an individual in the search space are represented by *chromosomes*, in essence a set of character strings. A *gene* is a subsection of a chromosome which encodes the value of a single parameter being optimized. Typical encodings for a gene could be *binary* or *integer*.

Through simulation of the evolutionary operations *recombination*, *mutation*, and *selection* new generations of search points are found that show a higher average fitness than their ancestors. [Figure 61.1](#) illustrates these steps.

Figure 61.1. Structure of a Genetic Algorithm

According to the comp.ai.genetic FAQ it cannot be stressed too strongly that a GA is not a pure random search for a solution to a problem. A GA uses stochastic processes, but the result is distinctly non-random (better than random).

61.3. Genetic Query Optimization (GEQO) in Postgres Pro

The GEQO module approaches the query optimization problem as though it were the well-known traveling salesman problem (TSP). Possible query plans are encoded as integer strings. Each string represents the join order from one relation of the query to the next. For example, the join tree

```

  /\
 /\ 2
 /\ 3

```


4 1

is encoded by the integer string '4-1-3-2', which means, first join relation '4' and '1', then '3', and then '2', where 1, 2, 3, 4 are relation IDs within the Postgres Pro optimizer.

Specific characteristics of the GEQO implementation in Postgres Pro are:

- Usage of a *steady state* GA (replacement of the least fit individuals in a population, not whole-generational replacement) allows fast convergence towards improved query plans. This is essential for query handling with reasonable time;
- Usage of *edge recombination crossover* which is especially suited to keep edge losses low for the solution of the TSP by means of a GA;
- Mutation as genetic operator is deprecated so that no repair mechanisms are needed to generate legal TSP tours.

Parts of the GEQO module are adapted from D. Whitley's Genitor algorithm.

The GEQO module allows the Postgres Pro query optimizer to support large join queries effectively through non-exhaustive search.

61.3.1. Generating Possible Plans with GEQO

The GEQO planning process uses the standard planner code to generate plans for scans of individual relations. Then join plans are developed using the genetic approach. As shown above, each candidate join plan is represented by a sequence in which to join the base relations. In the initial stage, the GEQO code simply generates some possible join sequences at random. For each join sequence considered, the standard planner code is invoked to estimate the cost of performing the query using that join sequence. (For each step of the join sequence, all three possible join strategies are considered; and all the initially-determined relation scan plans are available. The estimated cost is the cheapest of these possibilities.) Join sequences with lower estimated cost are considered “more fit” than those with higher cost. The genetic algorithm discards the least fit candidates. Then new candidates are generated by combining genes of more-fit candidates — that is, by using randomly-chosen portions of known low-cost join sequences to create new sequences for consideration. This process is repeated until a preset number of join sequences have been considered; then the best one found at any time during the search is used to generate the finished plan.

This process is inherently nondeterministic, because of the randomized choices made during both the initial population selection and subsequent “mutation” of the best candidates. To avoid surprising changes of the selected plan, each run of the GEQO algorithm restarts its random number generator with the current `geqo_seed` parameter setting. As long as `geqo_seed` and the other GEQO parameters are kept fixed, the same plan will be generated for a given query (and other planner inputs such as statistics). To experiment with different search paths, try changing `geqo_seed`.

61.3.2. Future Implementation Tasks for PostgreSQL GEQO

Work is still needed to improve the genetic algorithm parameter settings. We have to find a compromise for the parameter settings to satisfy two competing demands:

- Optimality of the query plan
- Computing time

In the current implementation, the fitness of each candidate join sequence is estimated by running the standard planner's join selection and cost estimation code from scratch. To the extent that different candidates use similar sub-sequences of joins, a great deal of work will be repeated. This could be made significantly faster by retaining cost estimates for sub-joins. The problem is to avoid expending unreasonable amounts of memory on retaining that state.

At a more basic level, it is not clear that solving query optimization with a GA algorithm designed for TSP is appropriate. In the TSP case, the cost associated with any substring (partial tour) is independent of the rest of the tour, but this is certainly not true for query optimization. Thus it is questionable whether edge recombination crossover is the most effective mutation procedure.

61.4. Further Reading

The following resources contain additional information about genetic algorithms:

- *The Hitch-Hiker's Guide to Evolutionary Computation*, (FAQ for <news://comp.ai.genetic>)
- *Evolutionary Computation and its application to art and design*, by Craig Reynolds
- [elma04](#)
- [fong](#)

Chapter 62. Table Access Method Interface Definition

This chapter explains the interface between the core Postgres Pro system and *table access methods*, which manage the storage for tables. The core system knows little about these access methods beyond what is specified here, so it is possible to develop entirely new access method types by writing add-on code.

Each table access method is described by a row in the `pg_am` system catalog. The `pg_am` entry specifies a name and a *handler function* for the table access method. These entries can be created and deleted using the [CREATE ACCESS METHOD](#) and [DROP ACCESS METHOD](#) SQL commands.

A table access method handler function must be declared to accept a single argument of type `internal` and to return the pseudo-type `table_am_handler`. The argument is a dummy value that simply serves to prevent handler functions from being called directly from SQL commands. The result of the function must be a pointer to a struct of type `TableAmRoutine`, which contains everything that the core code needs to know to make use of the table access method. The return value needs to be of server lifetime, which is typically achieved by defining it as a `static const` variable in global scope. The `TableAmRoutine` struct, also called the access method's *API struct*, defines the behavior of the access method using callbacks. These callbacks are pointers to plain C functions and are not visible or callable at the SQL level. All the callbacks and their behavior is defined in the `TableAmRoutine` structure (with comments inside the struct defining the requirements for callbacks). Most callbacks have wrapper functions, which are documented from the point of view of a user (rather than an implementor) of the table access method. For details, please refer to the [src/include/access/tableam.h](#) file.

To implement an access method, an implementor will typically need to implement an AM-specific type of tuple table slot (see [src/include/executor/tuptable.h](#)), which allows code outside the access method to hold references to tuples of the AM, and to access the columns of the tuple.

Currently, the way an AM actually stores data is fairly unconstrained. For example, it's possible, but not required, to use postgres' shared buffer cache. In case it is used, it likely makes sense to use Postgres Pro's standard page layout as described in [Section 66.6](#).

One fairly large constraint of the table access method API is that, currently, if the AM wants to support modifications and/or indexes, it is necessary for each tuple to have a tuple identifier (TID) consisting of a block number and an item number (see also [Section 66.6](#)). It is not strictly necessary that the sub-parts of TIDs have the same meaning they e.g., have for `heap`, but if bitmap scan support is desired (it is optional), the block number needs to provide locality.

For crash safety, an AM can use postgres' [WAL](#), or a custom implementation. If WAL is chosen, either [Generic WAL Records](#) can be used, or a [Custom WAL Resource Manager](#) can be implemented.

To implement transactional support in a manner that allows different table access methods be accessed within a single transaction, it likely is necessary to closely integrate with the machinery in [src/backend/access/transam/xlog.c](#).

Chapter 63. Index Access Method Interface Definition

This chapter defines the interface between the core Postgres Pro system and *index access methods*, which manage individual index types. The core system knows nothing about indexes beyond what is specified here, so it is possible to develop entirely new index types by writing add-on code.

All indexes in Postgres Pro are what are known technically as *secondary indexes*; that is, the index is physically separate from the table file that it describes. Each index is stored as its own physical *relation* and so is described by an entry in the `pg_class` catalog. The contents of an index are entirely under the control of its index access method. In practice, all index access methods divide indexes into standard-size pages so that they can use the regular storage manager and buffer manager to access the index contents. (All the existing index access methods furthermore use the standard page layout described in [Section 66.6](#), and most use the same format for index tuple headers; but these decisions are not forced on an access method.)

An index is effectively a mapping from some data key values to *tuple identifiers*, or TIDs, of row versions (tuples) in the index's parent table. A TID consists of a block number and an item number within that block (see [Section 66.6](#)). This is sufficient information to fetch a particular row version from the table. Indexes are not directly aware that under MVCC, there might be multiple extant versions of the same logical row; to an index, each tuple is an independent object that needs its own index entry. Thus, an update of a row always creates all-new index entries for the row, even if the key values did not change. ([HOT tuples](#) are an exception to this statement; but indexes do not deal with those, either.) Index entries for dead tuples are reclaimed (by vacuuming) when the dead tuples themselves are reclaimed.

63.1. Basic API Structure for Indexes

Each index access method is described by a row in the `pg_am` system catalog. The `pg_am` entry specifies a name and a *handler function* for the index access method. These entries can be created and deleted using the [CREATE ACCESS METHOD](#) and [DROP ACCESS METHOD](#) SQL commands.

An index access method handler function must be declared to accept a single argument of type `internal` and to return the pseudo-type `index_am_handler`. The argument is a dummy value that simply serves to prevent handler functions from being called directly from SQL commands. The result of the function must be a palloc'd struct of type `IndexAmRoutine`, which contains everything that the core code needs to know to make use of the index access method. The `IndexAmRoutine` struct, also called the access method's *API struct*, includes fields specifying assorted fixed properties of the access method, such as whether it can support multicolumn indexes. More importantly, it contains pointers to support functions for the access method, which do all of the real work to access indexes. These support functions are plain C functions and are not visible or callable at the SQL level. The support functions are described in [Section 63.2](#).

The structure `IndexAmRoutine` is defined thus:

```
typedef struct IndexAmRoutine
{
    NodeTag      type;

    /*
     * Total number of strategies (operators) by which we can traverse/search
     * this AM. Zero if AM does not have a fixed set of strategy assignments.
     */
    uint16       amstrategies;
    /* total number of support functions that this AM uses */
    uint16       amsupport;
    /* opclass options support function number or 0 */
    uint16       amoptsprocnum;
    /* does AM support ORDER BY indexed column's value? */

```

Index Access Method Interface Definition

```
bool            amcanorder;
/* does AM support ORDER BY result of an operator on indexed column? */
bool            amcanorderbyop;
/*
 * Does AM support only one ORDER BY operator on first indexed column?
 * amcanorderbyop is implied.
 */
bool            amorderbyopfirstcol;
/* does AM support backward scanning? */
bool            amcanbackward;
/* does AM support UNIQUE indexes? */
bool            amcanunique;
/* does AM support multi-column indexes? */
bool            amcanmulticol;
/* does AM require scans to have a constraint on the first index column? */
bool            amoptionalkey;
/* does AM handle ScalarArrayOpExpr quals? */
bool            amsearcharray;
/* does AM handle IS NULL/IS NOT NULL quals? */
bool            amsearchnulls;
/* can index storage data type differ from column data type? */
bool            amstorage;
/* can an index of this type be clustered on? */
bool            amclusterable;
/* does AM handle predicate locks? */
bool            ampredlocks;
/* does AM support parallel scan? */
bool            amcanparallel;
/* does AM support parallel build? */
bool            amcanbuildparallel;
/* does AM support columns included with clause INCLUDE? */
bool            amcaninclude;
/* does AM use maintenance_work_mem? */
bool            amusemaintenanceworkmem;
/* does AM summarize tuples, with at least all tuples in the block
 * summarized in one summary */
bool            amsummarizing;
/* OR of parallel vacuum flags */
uint8           amparallelvacuumoptions;
/* type of data stored in index, or InvalidOid if variable */
Oid             amkeytype;

/* interface functions */
ambuild_function ambuild;
ambuildempty_function ambuildempty;
aminert_function aminert;
aminertcleanup_function aminertcleanup;
ambulkdelete_function ambulkdelete;
amvacuumcleanup_function amvacuumcleanup;
amcanreturn_function amcanreturn; /* can be NULL */
amcostestimate_function amcostestimate;
amoptions_function amoptions;
amproperty_function amproperty; /* can be NULL */
ambuildphasename_function ambuildphasename; /* can be NULL */
amvalidate_function amvalidate;
amadjustmembers_function amadjustmembers; /* can be NULL */
ambeginscan_function ambeginscan;
amrescan_function amrescan;
```

Index Access Method Interface Definition

```
amgettupple_function amgettupple;      /* can be NULL */
amgetbitmap_function amgetbitmap;      /* can be NULL */
amendscan_function amendscan;
ammarkpos_function ammarkpos;          /* can be NULL */
amrestrpos_function amrestrpos;         /* can be NULL */

/* interface functions to support parallel index scans */
amestimateparallelscan_function amestimateparallelscan; /* can be NULL */
aminitparallelscan_function aminitparallelscan;         /* can be NULL */
amparallelrescan_function amparallelrescan;             /* can be NULL */
} IndexAmRoutine;
```

To be useful, an index access method must also have one or more *operator families* and *operator classes* defined in [pg_opfamily](#), [pg_opclass](#), [pg_amop](#), and [pg_amproc](#). These entries allow the planner to determine what kinds of query qualifications can be used with indexes of this access method. Operator families and classes are described in [Section 39.16](#), which is prerequisite material for reading this chapter.

An individual index is defined by a [pg_class](#) entry that describes it as a physical relation, plus a [pg_index](#) entry that shows the logical content of the index — that is, the set of index columns it has and the semantics of those columns, as captured by the associated operator classes. The index columns (key values) can be either simple columns of the underlying table or expressions over the table rows. The index access method normally has no interest in where the index key values come from (it is always handed precomputed key values) but it will be very interested in the operator class information in [pg_index](#). Both of these catalog entries can be accessed as part of the `Relation` data structure that is passed to all operations on the index.

Some of the flag fields of `IndexAmRoutine` have nonobvious implications. The requirements of `amcanunique` are discussed in [Section 63.5](#). The `amcanmulticol` flag asserts that the access method supports multi-key-column indexes, while `amoptionalkey` asserts that it allows scans where no indexable restriction clause is given for the first index column. When `amcanmulticol` is false, `amoptionalkey` essentially says whether the access method supports full-index scans without any restriction clause. Access methods that support multiple index columns *must* support scans that omit restrictions on any or all of the columns after the first; however they are permitted to require some restriction to appear for the first index column, and this is signaled by setting `amoptionalkey` false. One reason that an index AM might set `amoptionalkey` false is if it doesn't index null values. Since most indexable operators are strict and hence cannot return true for null inputs, it is at first sight attractive to not store index entries for null values: they could never be returned by an index scan anyway. However, this argument fails when an index scan has no restriction clause for a given index column. In practice this means that indexes that have `amoptionalkey` true must index nulls, since the planner might decide to use such an index with no scan keys at all. A related restriction is that an index access method that supports multiple index columns *must* support indexing null values in columns after the first, because the planner will assume the index can be used for queries that do not restrict these columns. For example, consider an index on (a,b) and a query with `WHERE a = 4`. The system will assume the index can be used to scan for rows with `a = 4`, which is wrong if the index omits rows where `b` is null. It is, however, OK to omit rows where the first indexed column is null. An index access method that does index nulls may also set `amsearchnulls`, indicating that it supports `IS NULL` and `IS NOT NULL` clauses as search conditions.

The `amcaninclude` flag indicates whether the access method supports “included” columns, that is it can store (without processing) additional columns beyond the key column(s). The requirements of the preceding paragraph apply only to the key columns. In particular, the combination of `amcanmulticol=false` and `amcaninclude=true` is sensible: it means that there can only be one key column, but there can also be included column(s). Also, included columns must be allowed to be null, independently of `amoptionalkey`.

The `amsummarizing` flag indicates whether the access method summarizes the indexed tuples, with summarizing granularity of at least per block. Access methods that do not point to individual tuples, but to block ranges (like BRIN), may allow the HOT optimization to continue. This does not apply to attributes referenced in index predicates, an update of such an attribute always disables HOT.

63.2. Index Access Method Functions

The index construction and maintenance functions that an index access method must provide in `IndexAmRoutine` are:

```
IndexBuildResult *
ambuild (Relation heapRelation,
         Relation indexRelation,
         IndexInfo *indexInfo);
```

Build a new index. The index relation has been physically created, but is empty. It must be filled in with whatever fixed data the access method requires, plus entries for all tuples already existing in the table. Ordinarily the `ambuild` function will call `table_index_build_scan()` to scan the table for existing tuples and compute the keys that need to be inserted into the index. The function must return a palloc'd struct containing statistics about the new index. The `amcanbuildparallel` flag indicates whether the access method supports parallel index builds. When set to `true`, the system will attempt to allocate parallel workers for the build. Access methods supporting only non-parallel index builds should leave this flag set to `false`.

```
void
ambuildempty (Relation indexRelation);
```

Build an empty index, and write it to the initialization fork (`INIT_FORKNUM`) of the given relation. This method is called only for unlogged indexes; the empty index written to the initialization fork will be copied over the main relation fork on each server restart.

```
bool
aminsert (Relation indexRelation,
          Datum *values,
          bool *isnull,
          ItemPointer heap_tid,
          Relation heapRelation,
          IndexUniqueCheck checkUnique,
          bool indexUnchanged,
          IndexInfo *indexInfo);
```

Insert a new tuple into an existing index. The `values` and `isnull` arrays give the key values to be indexed, and `heap_tid` is the TID to be indexed. If the access method supports unique indexes (its `amcanunique` flag is `true`) then `checkUnique` indicates the type of uniqueness check to perform. This varies depending on whether the unique constraint is deferrable; see [Section 63.5](#) for details. Normally the access method only needs the `heapRelation` parameter when performing uniqueness checking (since then it will have to look into the heap to verify tuple liveness).

The `indexUnchanged` Boolean value gives a hint about the nature of the tuple to be indexed. When it is `true`, the tuple is a duplicate of some existing tuple in the index. The new tuple is a logically unchanged successor MVCC tuple version. This happens when an `UPDATE` takes place that does not modify any columns covered by the index, but nevertheless requires a new version in the index. The index AM may use this hint to decide to apply bottom-up index deletion in parts of the index where many versions of the same logical row accumulate. Note that updating a non-key column or a column that only appears in a partial index predicate does not affect the value of `indexUnchanged`. The core code determines each tuple's `indexUnchanged` value using a low overhead approach that allows both false positives and false negatives. Index AMs must not treat `indexUnchanged` as an authoritative source of information about tuple visibility or versioning.

The function's Boolean result value is significant only when `checkUnique` is `UNIQUE_CHECK_PARTIAL`. In this case a `true` result means the new entry is known unique, whereas `false` means it might be non-unique (and a deferred uniqueness check must be scheduled). For other cases a constant `false` result is recommended.

Some indexes might not index all tuples. If the tuple is not to be indexed, `aminsert` should just return without doing anything.

If the index AM wishes to cache data across successive index insertions within an SQL statement, it can allocate space in `indexInfo->ii_Context` and store a pointer to the data in `indexInfo->ii_AmCache` (which will be NULL initially). If resources other than memory have to be released after index insertions, `aminertcleanup` may be provided, which will be called before the memory is released.

```
void
aminertcleanup (Relation indexRelation,
                IndexInfo *indexInfo);
```

Clean up state that was maintained across successive inserts in `indexInfo->ii_AmCache`. This is useful if the data requires additional cleanup steps (e.g., releasing pinned buffers), and simply releasing the memory is not sufficient.

```
IndexBulkDeleteResult *
ambulkdelete (IndexVacuumInfo *info,
              IndexBulkDeleteResult *stats,
              IndexBulkDeleteCallback callback,
              void *callback_state);
```

Delete tuple(s) from the index. This is a “bulk delete” operation that is intended to be implemented by scanning the whole index and checking each entry to see if it should be deleted. The passed-in `callback` function must be called, in the style `callback(TID, callback_state)` returns `bool`, to determine whether any particular index entry, as identified by its referenced TID, is to be deleted. Must return either NULL or a palloc'd struct containing statistics about the effects of the deletion operation. It is OK to return NULL if no information needs to be passed on to `amvacuumcleanup`.

Because of limited `maintenance_work_mem`, `ambulkdelete` might need to be called more than once when many tuples are to be deleted. The `stats` argument is the result of the previous call for this index (it is NULL for the first call within a `VACUUM` operation). This allows the AM to accumulate statistics across the whole operation. Typically, `ambulkdelete` will modify and return the same struct if the passed `stats` is not null.

```
IndexBulkDeleteResult *
amvacuumcleanup (IndexVacuumInfo *info,
                 IndexBulkDeleteResult *stats);
```

Clean up after a `VACUUM` operation (zero or more `ambulkdelete` calls). This does not have to do anything beyond returning index statistics, but it might perform bulk cleanup such as reclaiming empty index pages. `stats` is whatever the last `ambulkdelete` call returned, or NULL if `ambulkdelete` was not called because no tuples needed to be deleted. If the result is not NULL it must be a palloc'd struct. The statistics it contains will be used to update `pg_class`, and will be reported by `VACUUM` if `VERBOSE` is given. It is OK to return NULL if the index was not changed at all during the `VACUUM` operation, but otherwise correct stats should be returned.

`amvacuumcleanup` will also be called at completion of an `ANALYZE` operation. In this case `stats` is always NULL and any return value will be ignored. This case can be distinguished by checking `info->analyze_only`. It is recommended that the access method do nothing except post-insert cleanup in such a call, and that only in an autovacuum worker process.

```
bool
amcanreturn (Relation indexRelation, int attno);
```

Check whether the index can support *index-only scans* on the given column, by returning the column's original indexed value. The attribute number is 1-based, i.e., the first column's `attno` is 1. Returns true if supported, else false. This function should always return true for included columns (if those are supported), since there's little point in an included column that can't be retrieved. If the access method does not support index-only scans at all, the `amcanreturn` field in its `IndexAmRoutine` struct can be set to NULL.

```
void
amcostestimate (PlannerInfo *root,
```



```
IndexPath *path,
double loop_count,
Cost *indexStartupCost,
Cost *indexTotalCost,
Selectivity *indexSelectivity,
double *indexCorrelation,
double *indexPages);
```

Estimate the costs of an index scan. This function is described fully in [Section 63.6](#), below.

```
bytea *
amoptions (ArrayType *reloptions,
          bool validate);
```

Parse and validate the reloptions array for an index. This is called only when a non-null reloptions array exists for the index. *reloptions* is a text array containing entries of the form *name=value*. The function should construct a bytea value, which will be copied into the *rd_options* field of the index's relcache entry. The data contents of the bytea value are open for the access method to define; most of the standard access methods use struct *StdRdOptions*. When *validate* is true, the function should report a suitable error message if any of the options are unrecognized or have invalid values; when *validate* is false, invalid entries should be silently ignored. (*validate* is false when loading options already stored in *pg_catalog*; an invalid entry could only be found if the access method has changed its rules for options, and in that case ignoring obsolete entries is appropriate.) It is OK to return NULL if default behavior is wanted.

```
bool
amproperty (Oid index_oid, int attno,
            IndexAMProperty prop, const char *propname,
            bool *res, bool *isnull);
```

The *amproperty* method allows index access methods to override the default behavior of *pg_index_column_has_property* and related functions. If the access method does not have any special behavior for index property inquiries, the *amproperty* field in its *IndexAmRoutine* struct can be set to NULL. Otherwise, the *amproperty* method will be called with *index_oid* and *attno* both zero for *pg_indexam_has_property* calls, or with *index_oid* valid and *attno* zero for *pg_index_has_property* calls, or with *index_oid* valid and *attno* greater than zero for *pg_index_column_has_property* calls. *prop* is an enum value identifying the property being tested, while *propname* is the original property name string. If the core code does not recognize the property name then *prop* is *AMPROP_UNKNOWN*. Access methods can define custom property names by checking *propname* for a match (use *pg_strcasecmp* to match, for consistency with the core code); for names known to the core code, it's better to inspect *prop*. If the *amproperty* method returns true then it has determined the property test result: it must set **res* to the Boolean value to return, or set **isnull* to true to return a NULL. (Both of the referenced variables are initialized to false before the call.) If the *amproperty* method returns false then the core code will proceed with its normal logic for determining the property test result.

Access methods that support ordering operators should implement *AMPROP_DISTANCE_ORDERABLE* property testing, as the core code does not know how to do that and will return NULL. It may also be advantageous to implement *AMPROP_RETURNABLE* testing, if that can be done more cheaply than by opening the index and calling *amcanreturn*, which is the core code's default behavior. The default behavior should be satisfactory for all other standard properties.

```
char *
ambuildphasename (int64 phasenum);
```

Return the textual name of the given build phase number. The phase numbers are those reported during an index build via the *pgstat_progress_update_param* interface. The phase names are then exposed in the *pg_stat_progress_create_index* view.

```
bool
amvalidate (Oid opclassoid);
```

Validate the catalog entries for the specified operator class, so far as the access method can reasonably do that. For example, this might include testing that all required support functions are provided. The `amvalidate` function must return false if the `opclass` is invalid. Problems should be reported with `ereport` messages, typically at `INFO` level.

```
void
amadjustmembers (Oid opfamilyoid,
                 Oid opclassoid,
                 List *operators,
                 List *functions);
```

Validate proposed new operator and function members of an operator family, so far as the access method can reasonably do that, and set their dependency types if the default is not satisfactory. This is called during `CREATE OPERATOR CLASS` and during `ALTER OPERATOR FAMILY ADD`; in the latter case `opclassoid` is `InvalidOid`. The `List` arguments are lists of `OpFamilyMember` structs. Tests done by this function will typically be a subset of those performed by `amvalidate`, since `amadjustmembers` cannot assume that it is seeing a complete set of members. For example, it would be reasonable to check the signature of a support function, but not to check whether all required support functions are provided. Any problems can be reported by throwing an error. The dependency-related fields of the `OpFamilyMember` structs are initialized by the core code to create hard dependencies on the `opclass` if this is `CREATE OPERATOR CLASS`, or soft dependencies on the `opfamily` if this is `ALTER OPERATOR FAMILY ADD`. `amadjustmembers` can adjust these fields if some other behavior is more appropriate. For example, `GIN`, `GiST`, and `SP-GiST` always set operator members to have soft dependencies on the `opfamily`, since the connection between an operator and an `opclass` is relatively weak in these index types; so it is reasonable to allow operator members to be added and removed freely. Optional support functions are typically also given soft dependencies, so that they can be removed if necessary.

The purpose of an index, of course, is to support scans for tuples matching an indexable `WHERE` condition, often called a *qualifier* or *scan key*. The semantics of index scanning are described more fully in [Section 63.3](#), below. An index access method can support “plain” index scans, “bitmap” index scans, or both. The scan-related functions that an index access method must or may provide are:

```
IndexScanDesc
ambeginscan (Relation indexRelation,
            int nkeys,
            int norderbys);
```

Prepare for an index scan. The `nkeys` and `norderbys` parameters indicate the number of quals and ordering operators that will be used in the scan; these may be useful for space allocation purposes. Note that the actual values of the scan keys aren't provided yet. The result must be a palloc'd struct. For implementation reasons the index access method *must* create this struct by calling `RelationGetIndexScan()`. In most cases `ambeginscan` does little beyond making that call and perhaps acquiring locks; the interesting parts of index-scan startup are in `amrescan`.

```
void
amrescan (IndexScanDesc scan,
         ScanKey keys,
         int nkeys,
         ScanKey orderbys,
         int norderbys);
```

Start or restart an index scan, possibly with new scan keys. (To restart using previously-passed keys, `NULL` is passed for `keys` and/or `orderbys`.) Note that it is not allowed for the number of keys or order-by operators to be larger than what was passed to `ambeginscan`. In practice the restart feature is used when a new outer tuple is selected by a nested-loop join and so a new key comparison value is needed, but the scan key structure remains the same.

```
bool
amgettupletuple (IndexScanDesc scan,
                 ScanDirection direction);
```

Fetch the next tuple in the given scan, moving in the given direction (forward or backward in the index). Returns true if a tuple was obtained, false if no matching tuples remain. In the true case the tuple TID is stored into the `scan` structure. Note that “success” means only that the index contains an entry that matches the scan keys, not that the tuple necessarily still exists in the heap or will pass the caller's snapshot test. On success, `amgettupple` must also set `scan->xs_recheck` to true or false. False means it is certain that the index entry matches the scan keys. True means this is not certain, and the conditions represented by the scan keys must be rechecked against the heap tuple after fetching it. This provision supports “lossy” index operators. Note that rechecking will extend only to the scan conditions; a partial index predicate (if any) is never rechecked by `amgettupple` callers.

If the index supports [index-only scans](#) (i.e., `amcanreturn` returns true for any of its columns), then on success the AM must also check `scan->xs_want_itup`, and if that is true it must return the originally indexed data for the index entry. Columns for which `amcanreturn` returns false can be returned as nulls. The data can be returned in the form of an `IndexTuple` pointer stored at `scan->xs_itup`, with tuple descriptor `scan->xs_itupdesc`; or in the form of a `HeapTuple` pointer stored at `scan->xs_hitup`, with tuple descriptor `scan->xs_hitupdesc`. (The latter format should be used when reconstructing data that might possibly not fit into an `IndexTuple`.) In either case, management of the data referenced by the pointer is the access method's responsibility. The data must remain good at least until the next `amgettupple`, `amrescan`, or `amendscan` call for the scan.

The `amgettupple` function need only be provided if the access method supports “plain” index scans. If it doesn't, the `amgettupple` field in its `IndexAmRoutine` struct must be set to NULL.

```
int64
amgetbitmap (IndexScanDesc scan,
             TIDBitmap *tbm);
```

Fetch all tuples in the given scan and add them to the caller-supplied `TIDBitmap` (that is, OR the set of tuple IDs into whatever set is already in the bitmap). The number of tuples fetched is returned (this might be just an approximate count, for instance some AMs do not detect duplicates). While inserting tuple IDs into the bitmap, `amgetbitmap` can indicate that rechecking of the scan conditions is required for specific tuple IDs. This is analogous to the `xs_recheck` output parameter of `amgettupple`. Note: in the current implementation, support for this feature is conflated with support for lossy storage of the bitmap itself, and therefore callers recheck both the scan conditions and the partial index predicate (if any) for recheckable tuples. That might not always be true, however. `amgetbitmap` and `amgettupple` cannot be used in the same index scan; there are other restrictions too when using `amgetbitmap`, as explained in [Section 63.3](#).

The `amgetbitmap` function need only be provided if the access method supports “bitmap” index scans. If it doesn't, the `amgetbitmap` field in its `IndexAmRoutine` struct must be set to NULL.

```
void
amendscan (IndexScanDesc scan);
```

End a scan and release resources. The `scan` struct itself should not be freed, but any locks or pins taken internally by the access method must be released, as well as any other memory allocated by `ambeginscan` and other scan-related functions.

```
void
ammarkpos (IndexScanDesc scan);
```

Mark current scan position. The access method need only support one remembered scan position per scan.

The `ammarkpos` function need only be provided if the access method supports ordered scans. If it doesn't, the `ammarkpos` field in its `IndexAmRoutine` struct may be set to NULL.

```
void
amrestrpos (IndexScanDesc scan);
```

Restore the scan to the most recently marked position.

The `amrestrpos` function need only be provided if the access method supports ordered scans. If it doesn't, the `amrestrpos` field in its `IndexAmRoutine` struct may be set to `NULL`.

In addition to supporting ordinary index scans, some types of index may wish to support *parallel index scans*, which allow multiple backends to cooperate in performing an index scan. The index access method should arrange things so that each cooperating process returns a subset of the tuples that would be performed by an ordinary, non-parallel index scan, but in such a way that the union of those subsets is equal to the set of tuples that would be returned by an ordinary, non-parallel index scan. Furthermore, while there need not be any global ordering of tuples returned by a parallel scan, the ordering of that subset of tuples returned within each cooperating backend must match the requested ordering. The following functions may be implemented to support parallel index scans:

```
Size
amestimateparallelscale (int nkeys,
                        int norderbys);
```

Estimate and return the number of bytes of dynamic shared memory which the access method will be needed to perform a parallel scan. (This number is in addition to, not in lieu of, the amount of space needed for AM-independent data in `ParallelIndexScanDescData`.)

The `nkeys` and `norderbys` parameters indicate the number of quals and ordering operators that will be used in the scan; the same values will be passed to `amrescan`. Note that the actual values of the scan keys aren't provided yet.

It is not necessary to implement this function for access methods which do not support parallel scans or for which the number of additional bytes of storage required is zero.

```
void
aminitparallelscale (void *target);
```

This function will be called to initialize dynamic shared memory at the beginning of a parallel scan. `target` will point to at least the number of bytes previously returned by `amestimateparallelscale`, and this function may use that amount of space to store whatever data it wishes.

It is not necessary to implement this function for access methods which do not support parallel scans or in cases where the shared memory space required needs no initialization.

```
void
amparallelrescale (IndexScanDesc scan);
```

This function, if implemented, will be called when a parallel index scan must be restarted. It should reset any shared state set up by `aminitparallelscale` such that the scan will be restarted from the beginning.

63.3. Index Scanning

In an index scan, the index access method is responsible for regurgitating the TIDs of all the tuples it has been told about that match the *scan keys*. The access method is *not* involved in actually fetching those tuples from the index's parent table, nor in determining whether they pass the scan's visibility test or other conditions.

A scan key is the internal representation of a `WHERE` clause of the form *index_key operator constant*, where the index key is one of the columns of the index and the operator is one of the members of the operator family associated with that index column. An index scan has zero or more scan keys, which are implicitly ANDed — the returned tuples are expected to satisfy all the indicated conditions.

The access method can report that the index is *lossy*, or requires rechecks, for a particular query. This implies that the index scan will return all the entries that pass the scan key, plus possibly additional entries that do not. The core system's index-scan machinery will then apply the index conditions again to the heap tuple to verify whether or not it really should be selected. If the recheck option is not specified, the index scan must return exactly the set of matching entries.

Note that it is entirely up to the access method to ensure that it correctly finds all and only the entries passing all the given scan keys. Also, the core system will simply hand off all the `WHERE` clauses that match the index keys and operator families, without any semantic analysis to determine whether they are redundant or contradictory. As an example, given `WHERE x > 4 AND x > 14` where `x` is a b-tree indexed column, it is left to the b-tree `amrescan` function to realize that the first scan key is redundant and can be discarded. The extent of preprocessing needed during `amrescan` will depend on the extent to which the index access method needs to reduce the scan keys to a “normalized” form.

Some access methods return index entries in a well-defined order, others do not. There are actually two different ways that an access method can support sorted output:

- Access methods that always return entries in the natural ordering of their data (such as btree) should set `amcanorder` to true. Currently, such access methods must use btree-compatible strategy numbers for their equality and ordering operators.
- Access methods that support ordering operators should set `amcanorderbyop` to true. This indicates that the index is capable of returning entries in an order satisfying `ORDER BY index_key operator constant`. Scan modifiers of that form can be passed to `amrescan` as described previously. If the access method supports only one `ORDER BY` operator on the first indexed column, then it should set `amorderbyopfirstcol` to true.

The `amgettupple` function has a `direction` argument, which can be either `ForwardScanDirection` (the normal case) or `BackwardScanDirection`. If the first call after `amrescan` specifies `BackwardScanDirection`, then the set of matching index entries is to be scanned back-to-front rather than in the normal front-to-back direction, so `amgettupple` must return the last matching tuple in the index, rather than the first one as it normally would. (This will only occur for access methods that set `amcanorder` to true.) After the first call, `amgettupple` must be prepared to advance the scan in either direction from the most recently returned entry. (But if `amcanbackward` is false, all subsequent calls will have the same direction as the first one.)

Access methods that support ordered scans must support “marking” a position in a scan and later returning to the marked position. The same position might be restored multiple times. However, only one position need be remembered per scan; a new `ammarkpos` call overrides the previously marked position. An access method that does not support ordered scans need not provide `ammarkpos` and `amrestrpos` functions in `IndexAmRoutine`; set those pointers to NULL instead.

Both the scan position and the mark position (if any) must be maintained consistently in the face of concurrent insertions or deletions in the index. It is OK if a freshly-inserted entry is not returned by a scan that would have found the entry if it had existed when the scan started, or for the scan to return such an entry upon rescanning or backing up even though it had not been returned the first time through. Similarly, a concurrent delete might or might not be reflected in the results of a scan. What is important is that insertions or deletions not cause the scan to miss or multiply return entries that were not themselves being inserted or deleted.

If the index stores the original indexed data values (and not some lossy representation of them), it is useful to support [index-only scans](#), in which the index returns the actual data not just the TID of the heap tuple. This will only avoid I/O if the visibility map shows that the TID is on an all-visible page; else the heap tuple must be visited anyway to check MVCC visibility. But that is no concern of the access method's.

Instead of using `amgettupple`, an index scan can be done with `amgetbitmap` to fetch all tuples in one call. This can be noticeably more efficient than `amgettupple` because it allows avoiding lock/unlock cycles within the access method. In principle `amgetbitmap` should have the same effects as repeated `amgettupple` calls, but we impose several restrictions to simplify matters. First of all, `amgetbitmap` returns all tuples at once and marking or restoring scan positions isn't supported. Secondly, the tuples are returned in a bitmap which doesn't have any specific ordering, which is why `amgetbitmap` doesn't take a `direction` argument. (Ordering operators will never be supplied for such a scan, either.) Also, there is no provision for index-only scans with `amgetbitmap`, since there is no way to return the contents of index tuples. Finally, `amgetbitmap` does not guarantee any locking of the returned tuples, with implications spelled out in [Section 63.4](#).

Note that it is permitted for an access method to implement only `amgetbitmap` and not `amgettupple`, or vice versa, if its internal implementation is unsuited to one API or the other.

63.4. Index Locking Considerations

Index access methods must handle concurrent updates of the index by multiple processes. The core Postgres Pro system obtains `AccessShareLock` on the index during an index scan, and `RowExclusiveLock` when updating the index (including plain `VACUUM`). Since these lock types do not conflict, the access method is responsible for handling any fine-grained locking it might need. An `ACCESS EXCLUSIVE` lock on the index as a whole will be taken only during index creation, destruction, or `REINDEX (SHARE UPDATE EXCLUSIVE)` is taken instead with `CONCURRENTLY`).

Building an index type that supports concurrent updates usually requires extensive and subtle analysis of the required behavior.

Aside from the index's own internal consistency requirements, concurrent updates create issues about consistency between the parent table (the *heap*) and the index. Because Postgres Pro separates accesses and updates of the heap from those of the index, there are windows in which the index might be inconsistent with the heap. We handle this problem with the following rules:

- A new heap entry is made before making its index entries. (Therefore a concurrent index scan is likely to fail to see the heap entry. This is okay because the index reader would be uninterested in an uncommitted row anyway. But see [Section 63.5](#).)
- When a heap entry is to be deleted (by `VACUUM`), all its index entries must be removed first.
- An index scan must maintain a pin on the index page holding the item last returned by `amgettupple`, and `ambulkdelete` cannot delete entries from pages that are pinned by other backends. The need for this rule is explained below.

Without the third rule, it is possible for an index reader to see an index entry just before it is removed by `VACUUM`, and then to arrive at the corresponding heap entry after that was removed by `VACUUM`. This creates no serious problems if that item number is still unused when the reader reaches it, since an empty item slot will be ignored by `heap_fetch()`. But what if a third backend has already re-used the item slot for something else? When using an MVCC-compliant snapshot, there is no problem because the new occupant of the slot is certain to be too new to pass the snapshot test. However, with a non-MVCC-compliant snapshot (such as `SnapshotAny`), it would be possible to accept and return a row that does not in fact match the scan keys. We could defend against this scenario by requiring the scan keys to be rechecked against the heap row in all cases, but that is too expensive. Instead, we use a pin on an index page as a proxy to indicate that the reader might still be “in flight” from the index entry to the matching heap entry. Making `ambulkdelete` block on such a pin ensures that `VACUUM` cannot delete the heap entry before the reader is done with it. This solution costs little in run time, and adds blocking overhead only in the rare cases where there actually is a conflict.

This solution requires that index scans be “synchronous”: we have to fetch each heap tuple immediately after scanning the corresponding index entry. This is expensive for a number of reasons. An “asynchronous” scan in which we collect many TIDs from the index, and only visit the heap tuples sometime later, requires much less index locking overhead and can allow a more efficient heap access pattern. Per the above analysis, we must use the synchronous approach for non-MVCC-compliant snapshots, but an asynchronous scan is workable for a query using an MVCC snapshot.

In an `amgetbitmap` index scan, the access method does not keep an index pin on any of the returned tuples. Therefore it is only safe to use such scans with MVCC-compliant snapshots.

When the `ampredlocks` flag is not set, any scan using that index access method within a serializable transaction will acquire a nonblocking predicate lock on the full index. This will generate a read-write conflict with the insert of any tuple into that index by a concurrent serializable transaction. If certain patterns of read-write conflicts are detected among a set of concurrent serializable transactions, one of those transactions may be canceled to protect data integrity. When the flag is set, it indicates that the index access method implements finer-grained predicate locking, which will tend to reduce the frequency of such transaction cancellations.

63.5. Index Uniqueness Checks

Postgres Pro enforces SQL uniqueness constraints using *unique indexes*, which are indexes that disallow multiple entries with identical keys. An access method that supports this feature sets `amcanunique` true. (At present, only b-tree supports it.) Columns listed in the `INCLUDE` clause are not considered when enforcing uniqueness.

Because of MVCC, it is always necessary to allow duplicate entries to exist physically in an index: the entries might refer to successive versions of a single logical row. The behavior we actually want to enforce is that no MVCC snapshot could include two rows with equal index keys. This breaks down into the following cases that must be checked when inserting a new row into a unique index:

- If a conflicting valid row has been deleted by the current transaction, it's okay. (In particular, since an `UPDATE` always deletes the old row version before inserting the new version, this will allow an `UPDATE` on a row without changing the key.)
- If a conflicting row has been inserted by an as-yet-uncommitted transaction, the would-be inserter must wait to see if that transaction commits. If it rolls back then there is no conflict. If it commits without deleting the conflicting row again, there is a uniqueness violation. (In practice we just wait for the other transaction to end and then redo the visibility check in toto.)
- Similarly, if a conflicting valid row has been deleted by an as-yet-uncommitted transaction, the would-be inserter must wait for that transaction to commit or abort, and then repeat the test.

Furthermore, immediately before reporting a uniqueness violation according to the above rules, the access method must recheck the liveness of the row being inserted. If it is committed dead then no violation should be reported. (This case cannot occur during the ordinary scenario of inserting a row that's just been created by the current transaction. It can happen during `CREATE UNIQUE INDEX CONCURRENTLY`, however.)

We require the index access method to apply these tests itself, which means that it must reach into the heap to check the commit status of any row that is shown to have a duplicate key according to the index contents. This is without a doubt ugly and non-modular, but it saves redundant work: if we did a separate probe then the index lookup for a conflicting row would be essentially repeated while finding the place to insert the new row's index entry. What's more, there is no obvious way to avoid race conditions unless the conflict check is an integral part of insertion of the new index entry.

If the unique constraint is deferrable, there is additional complexity: we need to be able to insert an index entry for a new row, but defer any uniqueness-violation error until end of statement or even later. To avoid unnecessary repeat searches of the index, the index access method should do a preliminary uniqueness check during the initial insertion. If this shows that there is definitely no conflicting live tuple, we are done. Otherwise, we schedule a recheck to occur when it is time to enforce the constraint. If, at the time of the recheck, both the inserted tuple and some other tuple with the same key are live, then the error must be reported. (Note that for this purpose, “live” actually means “any tuple in the index entry's HOT chain is live”.) To implement this, the `aminsert` function is passed a `checkUnique` parameter having one of the following values:

- `UNIQUE_CHECK_NO` indicates that no uniqueness checking should be done (this is not a unique index).
- `UNIQUE_CHECK_YES` indicates that this is a non-deferrable unique index, and the uniqueness check must be done immediately, as described above.
- `UNIQUE_CHECK_PARTIAL` indicates that the unique constraint is deferrable. Postgres Pro will use this mode to insert each row's index entry. The access method must allow duplicate entries into the index, and report any potential duplicates by returning false from `aminsert`. For each row for which false is returned, a deferred recheck will be scheduled.

The access method must identify any rows which might violate the unique constraint, but it is not an error for it to report false positives. This allows the check to be done without waiting for other transactions to finish; conflicts reported here are not treated as errors and will be rechecked later, by which time they may no longer be conflicts.

- `UNIQUE_CHECK_EXISTING` indicates that this is a deferred recheck of a row that was reported as a potential uniqueness violation. Although this is implemented by calling `amininsert`, the access method must *not* insert a new index entry in this case. The index entry is already present. Rather, the access method must check to see if there is another live index entry. If so, and if the target row is also still live, report error.

It is recommended that in a `UNIQUE_CHECK_EXISTING` call, the access method further verify that the target row actually does have an existing entry in the index, and report error if not. This is a good idea because the index tuple values passed to `amininsert` will have been recomputed. If the index definition involves functions that are not really immutable, we might be checking the wrong area of the index. Checking that the target row is found in the recheck verifies that we are scanning for the same tuple values as were used in the original insertion.

63.6. Index Cost Estimation Functions

The `amcostestimate` function is given information describing a possible index scan, including lists of WHERE and ORDER BY clauses that have been determined to be usable with the index. It must return estimates of the cost of accessing the index and the selectivity of the WHERE clauses (that is, the fraction of parent-table rows that will be retrieved during the index scan). For simple cases, nearly all the work of the cost estimator can be done by calling standard routines in the optimizer; the point of having an `amcostestimate` function is to allow index access methods to provide index-type-specific knowledge, in case it is possible to improve on the standard estimates.

Each `amcostestimate` function must have the signature:

```
void
amcostestimate (PlannerInfo *root,
                IndexPath *path,
                double loop_count,
                Cost *indexStartupCost,
                Cost *indexTotalCost,
                Selectivity *indexSelectivity,
                double *indexCorrelation,
                double *indexPages);
```

The first three parameters are inputs:

root

The planner's information about the query being processed.

path

The index access path being considered. All fields except cost and selectivity values are valid.

loop_count

The number of repetitions of the index scan that should be factored into the cost estimates. This will typically be greater than one when considering a parameterized scan for use in the inside of a nestloop join. Note that the cost estimates should still be for just one scan; a larger *loop_count* means that it may be appropriate to allow for some caching effects across multiple scans.

The last five parameters are pass-by-reference outputs:

**indexStartupCost*

Set to cost of index start-up processing

**indexTotalCost*

Set to total cost of index processing

**indexSelectivity*

Set to index selectivity

**indexCorrelation*

Set to correlation coefficient between index scan order and underlying table's order

**indexPages*

Set to number of index leaf pages

Note that cost estimate functions must be written in C, not in SQL or any available procedural language, because they must access internal data structures of the planner/optimizer.

The index access costs should be computed using the parameters: a sequential disk block fetch has cost `seq_page_cost`, a nonsequential fetch has cost `random_page_cost`, and the cost of processing one index row should usually be taken as `cpu_index_tuple_cost`. In addition, an appropriate multiple of `cpu_operator_cost` should be charged for any comparison operators invoked during index processing (especially evaluation of the indexquals themselves).

The “start-up cost” is the part of the total scan cost that must be expended before we can begin to fetch the first row. For most indexes this can be taken as zero, but an index type with a high start-up cost might want to set it nonzero.

The *indexSelectivity* should be set to the estimated fraction of the parent table rows that will be retrieved during the index scan. In the case of a lossy query, this will typically be higher than the fraction of rows that actually pass the given qual conditions.

The *indexCorrelation* should be set to the correlation (ranging between -1.0 and 1.0) between the index order and the table order. This is used to adjust the estimate for the cost of fetching rows from the parent table.

The *indexPages* should be set to the number of leaf pages. This is used to estimate the number of workers for parallel index scan.

When *loop_count* is greater than one, the returned numbers should be averages expected for any one scan of the index.

Cost Estimation

A typical cost estimator will proceed as follows:

1. Estimate and return the fraction of parent-table rows that will be visited based on the given qual conditions. In the absence of any index-type-specific knowledge, use the standard optimizer function `clauselist_selectivity()`:

```
*indexSelectivity = clauselist_selectivity(root, path->indexquals,
                                           path->indexinfo->rel->reloid,
                                           JOIN_INNER, NULL);
```

2. Estimate the number of index rows that will be visited during the scan. For many index types this is the same as *indexSelectivity* times the number of rows in the index, but it might be more. (Note that the index's size in pages and rows is available from the `path->indexinfo` struct.)
3. Estimate the number of index pages that will be retrieved during the scan. This might be just *indexSelectivity* times the index's size in pages.
4. Compute the index access cost. A generic estimator might do this:

```
/*
 * Our generic assumption is that the index pages will be read
 * sequentially, so they cost seq_page_cost each, not random_page_cost.
 * Also, we charge for evaluation of the indexquals at each index row.
 * All the costs are assumed to be paid incrementally during the scan.
 */
cost_qual_eval(&index_qual_cost, path->indexquals, root);
```

```
*indexStartupCost = index_qual_cost.startup;  
*indexTotalCost = seq_page_cost * numIndexPages +  
    (cpu_index_tuple_cost + index_qual_cost.per_tuple) * numIndexTuples;
```

However, the above does not account for amortization of index reads across repeated index scans.

5. Estimate the index correlation. For a simple ordered index on a single field, this can be retrieved from `pg_statistic`. If the correlation is not known, the conservative estimate is zero (no correlation).

Chapter 64. Write Ahead Logging for Extensions

Certain extensions, principally extensions that implement custom access methods, may need to perform write-ahead logging in order to ensure crash-safety. Postgres Pro provides two ways for extensions to achieve this goal.

First, extensions can choose to use [generic WAL](#), a special type of WAL record which describes changes to pages in a generic way. This method is simple to implement and does not require that an extension library be loaded in order to apply the records. However, generic WAL records will be ignored when performing logical decoding.

Second, extensions can choose to use a [custom resource manager](#). This method is more flexible, supports logical decoding, and can sometimes generate much smaller write-ahead log records than would be possible with generic WAL. However, it is more complex for an extension to implement.

64.1. Generic WAL Records

Although all built-in WAL-logged modules have their own types of WAL records, there is also a generic WAL record type, which describes changes to pages in a generic way.

Note

Generic WAL records are ignored during [Logical Decoding](#). If logical decoding is required for your extension, consider a Custom WAL Resource Manager.

The API for constructing generic WAL records is defined in `access/generic_xlog.h` and implemented in `access/transam/generic_xlog.c`.

To perform a WAL-logged data update using the generic WAL record facility, follow these steps:

1. `state = GenericXLogStart(relation)` — start construction of a generic WAL record for the given relation.
2. `page = GenericXLogRegisterBuffer(state, buffer, flags)` — register a buffer to be modified within the current generic WAL record. This function returns a pointer to a temporary copy of the buffer's page, where modifications should be made. (Do not modify the buffer's contents directly.) The third argument is a bit mask of flags applicable to the operation. Currently the only such flag is `GENERIC_XLOG_FULL_IMAGE`, which indicates that a full-page image rather than a delta update should be included in the WAL record. Typically this flag would be set if the page is new or has been rewritten completely. `GenericXLogRegisterBuffer` can be repeated if the WAL-logged action needs to modify multiple pages.
3. Apply modifications to the page images obtained in the previous step.
4. `GenericXLogFinish(state)` — apply the changes to the buffers and emit the generic WAL record.

WAL record construction can be canceled between any of the above steps by calling `GenericXLogAbort(state)`. This will discard all changes to the page image copies.

Please note the following points when using the generic WAL record facility:

- No direct modifications of buffers are allowed! All modifications must be done in copies acquired from `GenericXLogRegisterBuffer()`. In other words, code that makes generic WAL records should never call `BufferGetPage()` for itself. However, it remains the caller's responsibility to pin/unpin and lock/unlock the buffers at appropriate times. Exclusive lock must be held on each target buffer from before `GenericXLogRegisterBuffer()` until after `GenericXLogFinish()`.

- Registrations of buffers (step 2) and modifications of page images (step 3) can be mixed freely, i.e., both steps may be repeated in any sequence. Keep in mind that buffers should be registered in the same order in which locks are to be obtained on them during replay.
- The maximum number of buffers that can be registered for a generic WAL record is `MAX_GENERIC_XLOG_PAGES`. An error will be thrown if this limit is exceeded.
- Generic WAL assumes that the pages to be modified have standard layout, and in particular that there is no useful data between `pd_lower` and `pd_upper`.
- Since you are modifying copies of buffer pages, `GenericXLogStart()` does not start a critical section. Thus, you can safely do memory allocation, error throwing, etc. between `GenericXLogStart()` and `GenericXLogFinish()`. The only actual critical section is present inside `GenericXLogFinish()`. There is no need to worry about calling `GenericXLogAbort()` during an error exit, either.
- `GenericXLogFinish()` takes care of marking buffers dirty and setting their LSNs. You do not need to do this explicitly.
- For unlogged relations, everything works the same except that no actual WAL record is emitted. Thus, you typically do not need to do any explicit checks for unlogged relations.
- The generic WAL redo function will acquire exclusive locks to buffers in the same order as they were registered. After redoing all changes, the locks will be released in the same order.
- If `GENERIC_XLOG_FULL_IMAGE` is not specified for a registered buffer, the generic WAL record contains a delta between the old and the new page images. This delta is based on byte-by-byte comparison. This is not very compact for the case of moving data within a page, and might be improved in the future.

64.2. Custom WAL Resource Managers

This section explains the interface between the core Postgres Pro system and custom WAL resource managers, which enable extensions to integrate directly with the [WAL](#).

An extension, especially a [Table Access Method](#) or [Index Access Method](#), may need to use WAL for recovery, replication, and/or [Logical Decoding](#).

To create a new custom WAL resource manager, first define an `RmgrData` structure with implementations for the resource manager methods.

```
/*
 * Method table for resource managers.
 *
 * This struct must be kept in sync with the PG_RMGR definition in
 * rmgr.c.
 *
 * rm_identify must return a name for the record based on xl_info (without
 * reference to the rmid). For example, XLOG_BTREE_VACUUM would be named
 * "VACUUM". rm_desc can then be called to obtain additional detail for the
 * record, if available (e.g. the last block).
 *
 * rm_mask takes as input a page modified by the resource manager and masks
 * out bits that shouldn't be flagged by wal_consistency_checking.
 *
 * RmgrTable[] is indexed by RmgrId values (see rmgrlist.h). If rm_name is
 * NULL, the corresponding RmgrTable entry is considered invalid.
 */
typedef struct RmgrData
{
    const char *rm_name;
    void (*rm_redo) (XLogReaderState *record);
    void (*rm_desc) (StringInfo buf, XLogReaderState *record);
}
```

```
const char *(*rm_identify) (uint8 info);
void        (*rm_startup) (void);
void        (*rm_cleanup) (void);
void        (*rm_mask) (char *pagedata, BlockNumber blkno);
void        (*rm_decode) (struct LogicalDecodingContext *ctx,
                          struct XLogRecordBuffer *buf);

} RmgrData;
```

Then, register your new resource manager.

```
/*
 * Register a new custom WAL resource manager.
 *
 * Resource manager IDs must be globally unique across all extensions. Refer
 * to https://wiki.postgresql.org/wiki/CustomWALResourceManagers to reserve a
 * unique RmgrId for your extension, to avoid conflicts with other extension
 * developers. During development, use RM_EXPERIMENTAL_ID to avoid needlessly
 * reserving a new ID.
 */
extern void RegisterCustomRmgr(RmgrId rmid, const RmgrData *rmgr);
```

`RegisterCustomRmgr` must be called from the extension module's `_PG_init` function. While developing a new extension, use `RM_EXPERIMENTAL_ID` for `rmid`. When you are ready to release the extension to users, reserve a new resource manager ID at the [Custom WAL Resource Manager](#) page.

Place the extension module implementing the custom resource manager in [shared_preload_libraries](#) so that it will be loaded early during Postgres Pro startup.

Note

The extension must remain in `shared_preload_libraries` as long as any custom WAL records may exist in the system. Otherwise Postgres Pro will not be able to apply or decode the custom WAL records, which may prevent the server from starting.

Chapter 65. Built-in Index Access Methods

65.1. B-Tree Indexes

65.1.1. Introduction

Postgres Pro includes an implementation of the standard btree (multi-way balanced tree) index data structure. Any data type that can be sorted into a well-defined linear order can be indexed by a btree index. The only limitation is that an index entry cannot exceed approximately one-third of a page (after TOAST compression, if applicable).

Because each btree operator class imposes a sort order on its data type, btree operator classes (or, really, operator families) have come to be used as Postgres Pro's general representation and understanding of sorting semantics. Therefore, they've acquired some features that go beyond what would be needed just to support btree indexes, and parts of the system that are quite distant from the btree AM make use of them.

65.1.2. Behavior of B-Tree Operator Classes

As shown in [Table 39.3](#), a btree operator class must provide five comparison operators, `<`, `<=`, `=`, `>=` and `>`. One might expect that `<>` should also be part of the operator class, but it is not, because it would almost never be useful to use a `<>` WHERE clause in an index search. (For some purposes, the planner treats `<>` as associated with a btree operator class; but it finds that operator via the `=` operator's negator link, rather than from `pg_amop`.)

When several data types share near-identical sorting semantics, their operator classes can be grouped into an operator family. Doing so is advantageous because it allows the planner to make deductions about cross-type comparisons. Each operator class within the family should contain the single-type operators (and associated support functions) for its input data type, while cross-type comparison operators and support functions are “loose” in the family. It is recommendable that a complete set of cross-type operators be included in the family, thus ensuring that the planner can represent any comparison conditions that it deduces from transitivity.

There are some basic assumptions that a btree operator family must satisfy:

- An `=` operator must be an equivalence relation; that is, for all non-null values *A*, *B*, *C* of the data type:
 - *A* = *A* is true (*reflexive law*)
 - if *A* = *B*, then *B* = *A* (*symmetric law*)
 - if *A* = *B* and *B* = *C*, then *A* = *C* (*transitive law*)
- A `<` operator must be a strong ordering relation; that is, for all non-null values *A*, *B*, *C*:
 - *A* < *A* is false (*irreflexive law*)
 - if *A* < *B* and *B* < *C*, then *A* < *C* (*transitive law*)
- Furthermore, the ordering is total; that is, for all non-null values *A*, *B*:
 - exactly one of *A* < *B*, *A* = *B*, and *B* < *A* is true (*trichotomy law*)

(The trichotomy law justifies the definition of the comparison support function, of course.)

The other three operators are defined in terms of `=` and `<` in the obvious way, and must act consistently with them.

For an operator family supporting multiple data types, the above laws must hold when *A*, *B*, *C* are taken from any data types in the family. The transitive laws are the trickiest to ensure, as in cross-type situations they represent statements that the behaviors of two or three different operators are consistent. As

an example, it would not work to put `float8` and `numeric` into the same operator family, at least not with the current semantics that `numeric` values are converted to `float8` for comparison to a `float8`. Because of the limited accuracy of `float8`, this means there are distinct `numeric` values that will compare equal to the same `float8` value, and thus the transitive law would fail.

Another requirement for a multiple-data-type family is that any implicit or binary-coercion casts that are defined between data types included in the operator family must not change the associated sort ordering.

It should be fairly clear why a btree index requires these laws to hold within a single data type: without them there is no ordering to arrange the keys with. Also, index searches using a comparison key of a different data type require comparisons to behave sanely across two data types. The extensions to three or more data types within a family are not strictly required by the btree index mechanism itself, but the planner relies on them for optimization purposes.

To implement the distance ordered (nearest-neighbor) search, we only need to define a distance operator (usually called `<->`) with the corresponding operator family for distance comparison in the index's operator class. These operators must satisfy the following assumptions for all non-null values *A*, *B*, *C* of the datatype:

- $A <-> B = B <-> A$ (*symmetric law*)
- if $A = B$, then $A <-> C = B <-> C$ (*distance equivalence*)
- if $(A \leq B \text{ and } B \leq C)$ or $(A \geq B \text{ and } B \geq C)$, then $A <-> B \leq A <-> C$ (*monotonicity*)

65.1.3. B-Tree Support Functions

As shown in [Table 39.9](#), btree defines one required and four optional support functions. The five user-defined methods are:

`order`

For each combination of data types that a btree operator family provides comparison operators for, it must provide a comparison support function, registered in `pg_amproc` with support function number 1 and `amproclefttype/amprocrighttype` equal to the left and right data types for the comparison (i.e., the same data types that the matching operators are registered with in `pg_amop`). The comparison function must take two non-null values *A* and *B* and return an `int32` value that is `< 0`, `0`, or `> 0` when $A < B$, $A = B$, or $A > B$, respectively. A null result is disallowed: all values of the data type must be comparable.

If the compared values are of a collatable data type, the appropriate collation OID will be passed to the comparison support function, using the standard `PG_GET_COLLATION()` mechanism.

`sortsupport`

Optionally, a btree operator family may provide *sort support* function(s), registered under support function number 2. These functions allow implementing comparisons for sorting purposes in a more efficient way than naively calling the comparison support function. The APIs involved in this are defined in [src/include/utils/sortsupport.h](#).

`in_range`

Optionally, a btree operator family may provide *in_range* support function(s), registered under support function number 3. These are not used during btree index operations; rather, they extend the semantics of the operator family so that it can support window clauses containing the `RANGE offset PRECEDING` and `RANGE offset FOLLOWING` frame bound types (see [Section 4.2.8](#)). Fundamentally, the extra information provided is how to add or subtract an *offset* value in a way that is compatible with the family's data ordering.

An `in_range` function must have the signature

```
in_range(val type1, base type1, offset type2, sub bool, less bool)
```

returns bool

val and *base* must be of the same type, which is one of the types supported by the operator family (i.e., a type for which it provides an ordering). However, *offset* could be of a different type, which might be one otherwise unsupported by the family. An example is that the built-in `time_ops` family provides an `in_range` function that has *offset* of type `interval`. A family can provide `in_range` functions for any of its supported types and one or more *offset* types. Each `in_range` function should be entered in `pg_amproc` with `amproclefttype` equal to `type1` and `amprocrighttype` equal to `type2`.

The essential semantics of an `in_range` function depend on the two Boolean flag parameters. It should add or subtract *base* and *offset*, then compare *val* to the result, as follows:

- if *!sub* and *!less*, return $val \geq (base + offset)$
- if *!sub* and *less*, return $val \leq (base + offset)$
- if *sub* and *!less*, return $val \geq (base - offset)$
- if *sub* and *less*, return $val \leq (base - offset)$

Before doing so, the function should check the sign of *offset*: if it is less than zero, raise error `ERCODE_INVALID_PRECEDING_OR_FOLLOWING_SIZE` (22013) with error text like “invalid preceding or following size in window function”. (This is required by the SQL standard, although nonstandard operator families might perhaps choose to ignore this restriction, since there seems to be little semantic necessity for it.) This requirement is delegated to the `in_range` function so that the core code needn't understand what “less than zero” means for a particular data type.

An additional expectation is that `in_range` functions should, if practical, avoid throwing an error if $base + offset$ or $base - offset$ would overflow. The correct comparison result can be determined even if that value would be out of the data type's range. Note that if the data type includes concepts such as “infinity” or “NaN”, extra care may be needed to ensure that `in_range`'s results agree with the normal sort order of the operator family.

The results of the `in_range` function must be consistent with the sort ordering imposed by the operator family. To be precise, given any fixed values of *offset* and *sub*, then:

- If `in_range` with *less* = true is true for some *val1* and *base*, it must be true for every $val2 \leq val1$ with the same *base*.
- If `in_range` with *less* = true is false for some *val1* and *base*, it must be false for every $val2 \geq val1$ with the same *base*.
- If `in_range` with *less* = true is true for some *val* and *base1*, it must be true for every $base2 \geq base1$ with the same *val*.
- If `in_range` with *less* = true is false for some *val* and *base1*, it must be false for every $base2 \leq base1$ with the same *val*.

Analogous statements with inverted conditions hold when *less* = false.

If the type being ordered (*type1*) is collatable, the appropriate collation OID will be passed to the `in_range` function, using the standard `PG_GET_COLLATION()` mechanism.

`in_range` functions need not handle NULL inputs, and typically will be marked strict.

`equalimage`

Optionally, a btree operator family may provide `equalimage` (“equality implies image equality”) support functions, registered under support function number 4. These functions allow the core code to determine when it is safe to apply the btree deduplication optimization. Currently, `equalimage` functions are only called when building or rebuilding an index.

An `equalimage` function must have the signature

`equalimage(opcintype oid) returns bool`

The return value is static information about an operator class and collation. Returning `true` indicates that the `order` function for the operator class is guaranteed to only return 0 (“arguments are equal”) when its *A* and *B* arguments are also interchangeable without any loss of semantic information. Not registering an `equalimage` function or returning `false` indicates that this condition cannot be assumed to hold.

The `opcintype` argument is the `pg_type.oid` of the data type that the operator class indexes. This is a convenience that allows reuse of the same underlying `equalimage` function across operator classes. If `opcintype` is a collatable data type, the appropriate collation OID will be passed to the `equalimage` function, using the standard `PG_GET_COLLATION()` mechanism.

As far as the operator class is concerned, returning `true` indicates that deduplication is safe (or safe for the collation whose OID was passed to its `equalimage` function). However, the core code will only deem deduplication safe for an index when *every* indexed column uses an operator class that registers an `equalimage` function, and each function actually returns `true` when called.

Image equality is *almost* the same condition as simple bitwise equality. There is one subtle difference: When indexing a varlena data type, the on-disk representation of two image equal datums may not be bitwise equal due to inconsistent application of TOAST compression on input. Formally, when an operator class's `equalimage` function returns `true`, it is safe to assume that the `datum_image_eq()` C function will always agree with the operator class's `order` function (provided that the same collation OID is passed to both the `equalimage` and `order` functions).

The core code is fundamentally unable to deduce anything about the “equality implies image equality” status of an operator class within a multiple-data-type family based on details from other operator classes in the same family. Also, it is not sensible for an operator family to register a cross-type `equalimage` function, and attempting to do so will result in an error. This is because “equality implies image equality” status does not just depend on sorting/equality semantics, which are more or less defined at the operator family level. In general, the semantics that one particular data type implements must be considered separately.

The convention followed by the operator classes included with the core Postgres Pro distribution is to register a stock, generic `equalimage` function. Most operator classes register `btequalimage()`, which indicates that deduplication is safe unconditionally. Operator classes for collatable data types such as text register `btvarstrequalimage()`, which indicates that deduplication is safe with deterministic collations. Best practice for third-party extensions is to register their own custom function to retain control.

options

Optionally, a B-tree operator family may provide `options` (“operator class specific options”) support functions, registered under support function number 5. These functions define a set of user-visible parameters that control operator class behavior.

An `options` support function must have the signature

```
options(relopts local_relopts *) returns void
```

The function is passed a pointer to a `local_relopts` struct, which needs to be filled with a set of operator class specific options. The options can be accessed from other support functions using the `PG_HAS_OPCLASS_OPTIONS()` and `PG_GET_OPCLASS_OPTIONS()` macros.

Currently, no B-Tree operator class has an `options` support function. B-tree doesn't allow flexible representation of keys like GiST, SP-GiST, GIN and BRIN do. So, `options` probably doesn't have much application in the current B-tree index access method. Nevertheless, this support function was added to B-tree for uniformity, and will probably find uses during further evolution of B-tree in Postgres Pro.

65.1.4. Implementation

This section covers B-Tree index implementation details that may be of use to advanced users.

65.1.4.1. B-Tree Structure

Postgres Pro B-Tree indexes are multi-level tree structures, where each level of the tree can be used as a doubly-linked list of pages. A single metapage is stored in a fixed position at the start of the first segment file of the index. All other pages are either leaf pages or internal pages. Leaf pages are the pages on the lowest level of the tree. All other levels consist of internal pages. Each leaf page contains tuples that point to table rows. Each internal page contains tuples that point to the next level down in the tree. Typically, over 99% of all pages are leaf pages. Both internal pages and leaf pages use the standard page format described in [Section 66.6](#).

New leaf pages are added to a B-Tree index when an existing leaf page cannot fit an incoming tuple. A *page split* operation makes room for items that originally belonged on the overflowing page by moving a portion of the items to a new page. Page splits must also insert a new *downlink* to the new page in the parent page, which may cause the parent to split in turn. Page splits “cascade upwards” in a recursive fashion. When the root page finally cannot fit a new downlink, a *root page split* operation takes place. This adds a new level to the tree structure by creating a new root page that is one level above the original root page.

65.1.4.2. Bottom-up Index Deletion

B-Tree indexes are not directly aware that under MVCC, there might be multiple extant versions of the same logical table row; to an index, each tuple is an independent object that needs its own index entry. “Version churn” tuples may sometimes accumulate and adversely affect query latency and throughput. This typically occurs with `UPDATE`-heavy workloads where most individual updates cannot apply the [HOT optimization](#). Changing the value of only one column covered by one index during an `UPDATE` *always* necessitates a new set of index tuples — one for *each and every* index on the table. Note in particular that this includes indexes that were not “logically modified” by the `UPDATE`. All indexes will need a successor physical index tuple that points to the latest version in the table. Each new tuple within each index will generally need to coexist with the original “updated” tuple for a short period of time (typically until shortly after the `UPDATE` transaction commits).

B-Tree indexes incrementally delete version churn index tuples by performing *bottom-up index deletion* passes. Each deletion pass is triggered in reaction to an anticipated “version churn page split”. This only happens with indexes that are not logically modified by `UPDATE` statements, where concentrated build up of obsolete versions in particular pages would occur otherwise. A page split will usually be avoided, though it's possible that certain implementation-level heuristics will fail to identify and delete even one garbage index tuple (in which case a page split or deduplication pass resolves the issue of an incoming new tuple not fitting on a leaf page). The worst-case number of versions that any index scan must traverse (for any single logical row) is an important contributor to overall system responsiveness and throughput. A bottom-up index deletion pass targets suspected garbage tuples in a single leaf page based on *qualitative* distinctions involving logical rows and versions. This contrasts with the “top-down” index cleanup performed by autovacuum workers, which is triggered when certain *quantitative* table-level thresholds are exceeded (see [Section 24.1.6](#)).

Note

Not all deletion operations that are performed within B-Tree indexes are bottom-up deletion operations. There is a distinct category of index tuple deletion: *simple index tuple deletion*. This is a deferred maintenance operation that deletes index tuples that are known to be safe to delete (those whose item identifier's `LP_DEAD` bit is already set). Like bottom-up index deletion, simple index deletion takes place at the point that a page split is anticipated as a way of avoiding the split.

Simple deletion is opportunistic in the sense that it can only take place when recent index scans set the `LP_DEAD` bits of affected items in passing. Prior to Postgres Pro 14, the only category of B-Tree deletion was simple deletion. The main differences between it and bottom-up deletion are that only the former is opportunistically driven by the activity of passing index scans, while only the latter specifically targets version churn from `UPDATES` that do not logically modify indexed columns.

Bottom-up index deletion performs the vast majority of all garbage index tuple cleanup for particular indexes with certain workloads. This is expected with any B-Tree index that is subject to significant version churn from `UPDATE`s that rarely or never logically modify the columns that the index covers. The average and worst-case number of versions per logical row can be kept low purely through targeted incremental deletion passes. It's quite possible that the on-disk size of certain indexes will never increase by even one single page/block despite *constant* version churn from `UPDATE`s. Even then, an exhaustive “clean sweep” by a `VACUUM` operation (typically run in an autovacuum worker process) will eventually be required as a part of *collective* cleanup of the table and each of its indexes.

Unlike `VACUUM`, bottom-up index deletion does not provide any strong guarantees about how old the oldest garbage index tuple may be. No index can be permitted to retain “floating garbage” index tuples that became dead prior to a conservative cutoff point shared by the table and all of its indexes collectively. This fundamental table-level invariant makes it safe to recycle table TIDs. This is how it is possible for distinct logical rows to reuse the same table TID over time (though this can never happen with two logical rows whose lifetimes span the same `VACUUM` cycle).

65.1.4.3. Deduplication

A duplicate is a leaf page tuple (a tuple that points to a table row) where *all* indexed key columns have values that match corresponding column values from at least one other leaf page tuple in the same index. Duplicate tuples are quite common in practice. B-Tree indexes can use a special, space-efficient representation for duplicates when an optional technique is enabled: *deduplication*.

Deduplication works by periodically merging groups of duplicate tuples together, forming a single *posting list* tuple for each group. The column key value(s) only appear once in this representation. This is followed by a sorted array of TIDs that point to rows in the table. This significantly reduces the storage size of indexes where each value (or each distinct combination of column values) appears several times on average. The latency of queries can be reduced significantly. Overall query throughput may increase significantly. The overhead of routine index vacuuming may also be reduced significantly.

Note

B-Tree deduplication is just as effective with “duplicates” that contain a `NULL` value, even though `NULL` values are never equal to each other according to the `=` member of any B-Tree operator class. As far as any part of the implementation that understands the on-disk B-Tree structure is concerned, `NULL` is just another value from the domain of indexed values.

The deduplication process occurs lazily, when a new item is inserted that cannot fit on an existing leaf page, though only when index tuple deletion could not free sufficient space for the new item (typically deletion is briefly considered and then skipped over). Unlike GIN posting list tuples, B-Tree posting list tuples do not need to expand every time a new duplicate is inserted; they are merely an alternative physical representation of the original logical contents of the leaf page. This design prioritizes consistent performance with mixed read-write workloads. Most client applications will at least see a moderate performance benefit from using deduplication. Deduplication is enabled by default.

`CREATE INDEX` and `REINDEX` apply deduplication to create posting list tuples, though the strategy they use is slightly different. Each group of duplicate ordinary tuples encountered in the sorted input taken from the table is merged into a posting list tuple *before* being added to the current pending leaf page. Individual posting list tuples are packed with as many TIDs as possible. Leaf pages are written out in the usual way, without any separate deduplication pass. This strategy is well-suited to `CREATE INDEX` and `REINDEX` because they are once-off batch operations.

Write-heavy workloads that don't benefit from deduplication due to having few or no duplicate values in indexes will incur a small, fixed performance penalty (unless deduplication is explicitly disabled). The `deduplicate_items` storage parameter can be used to disable deduplication within individual indexes. There is never any performance penalty with read-only workloads, since reading posting list tuples is at least as efficient as reading the standard tuple representation. Disabling deduplication isn't usually helpful.

It is sometimes possible for unique indexes (as well as unique constraints) to use deduplication. This allows leaf pages to temporarily “absorb” extra version churn duplicates. Deduplication in unique indexes augments bottom-up index deletion, especially in cases where a long-running transaction holds a snapshot that blocks garbage collection. The goal is to buy time for the bottom-up index deletion strategy to become effective again. Delaying page splits until a single long-running transaction naturally goes away can allow a bottom-up deletion pass to succeed where an earlier deletion pass failed.

Tip

A special heuristic is applied to determine whether a deduplication pass in a unique index should take place. It can often skip straight to splitting a leaf page, avoiding a performance penalty from wasting cycles on unhelpful deduplication passes. If you're concerned about the overhead of deduplication, consider setting `deduplicate_items = off` selectively. Leaving deduplication enabled in unique indexes has little downside.

Deduplication cannot be used in all cases due to implementation-level restrictions. Deduplication safety is determined when `CREATE INDEX` or `REINDEX` is run.

Note that deduplication is deemed unsafe and cannot be used in the following cases involving semantically significant differences among equal datums:

- `text`, `varchar`, and `char` cannot use deduplication when a *nondeterministic* collation is used. Case and accent differences must be preserved among equal datums.
- `numeric` cannot use deduplication. Numeric display scale must be preserved among equal datums.
- `jsonb` cannot use deduplication, since the `jsonb` B-Tree operator class uses `numeric` internally.
- `float4` and `float8` cannot use deduplication. These types have distinct representations for `-0` and `0`, which are nevertheless considered equal. This difference must be preserved.

There is one further implementation-level restriction that may be lifted in a future version of Postgres Pro:

- Container types (such as composite types, arrays, or range types) cannot use deduplication.

There is one further implementation-level restriction that applies regardless of the operator class or collation used:

- `INCLUDE` indexes can never use deduplication.

65.2. GiST Indexes

65.2.1. Introduction

GiST stands for Generalized Search Tree. It is a balanced, tree-structured access method, that acts as a base template in which to implement arbitrary indexing schemes. B-trees, R-trees and many other indexing schemes can be implemented in GiST.

One advantage of GiST is that it allows the development of custom data types with the appropriate access methods, by an expert in the domain of the data type, rather than a database expert.

Some of the information here is derived from the University of California at Berkeley's GiST Indexing Project [web site](#) and Marcel Kornacker's thesis, [Access Methods for Next-Generation Database Systems](#). The GiST implementation in Postgres Pro is primarily maintained by Teodor Sigaev and Oleg Bartunov, and there is more information on their [web site](#).

65.2.2. Built-in Operator Classes

The core Postgres Pro distribution includes the GiST operator classes shown in [Table 65.1](#). (Some of the optional modules described in [Appendix F](#) provide additional GiST operator classes.)

Table 65.1. Built-in GiST Operator Classes

Name	Indexable Operators	Ordering Operators
box_ops	<< (box, box)	<-> (box, point)
	&< (box, box)	
	&& (box, box)	
	&> (box, box)	
	>> (box, box)	
	~= (box, box)	
	@> (box, box)	
	<@ (box, box)	
	&< (box, box)	
	<< (box, box)	
	>> (box, box)	
	&> (box, box)	
circle_ops	<< (circle, circle)	<-> (circle, point)
	&< (circle, circle)	
	&> (circle, circle)	
	>> (circle, circle)	
	<@ (circle, circle)	
	@> (circle, circle)	
	~= (circle, circle)	
	&& (circle, circle)	
	>> (circle, circle)	
	<< (circle, circle)	
	&< (circle, circle)	
	&> (circle, circle)	
inet_ops	<< (inet, inet)	
	<<= (inet, inet)	
	>> (inet, inet)	
	>>= (inet, inet)	
	= (inet, inet)	
	<> (inet, inet)	
	< (inet, inet)	
	<= (inet, inet)	
	> (inet, inet)	
	>= (inet, inet)	
	&& (inet, inet)	
multirange_ops	= (anymultirange, anymultirange)	
	&& (anymultirange, anymultirange)	
	&& (anymultirange, anyrange)	
	@> (anymultirange, anyelement)	

Name	Indexable Operators	Ordering Operators
	@> (anymultirange, anymultirange) @> (anymultirange, anyrange) <@ (anymultirange, anymultirange) <@ (anymultirange, anyrange) << (anymultirange, anymultirange) << (anymultirange, anyrange) >> (anymultirange, anymultirange) >> (anymultirange, anyrange) &< (anymultirange, anymultirange) &< (anymultirange, anyrange) &> (anymultirange, anymultirange) &> (anymultirange, anyrange) - - (anymultirange, anymultirange) - - (anymultirange, anyrange)	
point_ops	>> (point, point) << (point, point) >> (point, point) << (point, point) ~ = (point, point) <@ (point, box) <@ (point, polygon) <@ (point, circle)	<-> (point, point)
poly_ops	<< (polygon, polygon) &< (polygon, polygon) &> (polygon, polygon) >> (polygon, polygon) <@ (polygon, polygon) @> (polygon, polygon) ~ = (polygon, polygon) && (polygon, polygon) << (polygon, polygon) &< (polygon, polygon) &> (polygon, polygon) >> (polygon, polygon)	<-> (polygon, point)
range_ops	= (anyrange, anyrange) && (anyrange, anyrange) && (anyrange, anymultirange) @> (anyrange, anyelement) @> (anyrange, anyrange) @> (anyrange, anymultirange)	

Name	Indexable Operators	Ordering Operators
	<@ (anyrange, anyrange)	
	<@ (anyrange, anymultirange)	
	<< (anyrange, anyrange)	
	<< (anyrange, anymultirange)	
	>> (anyrange, anyrange)	
	>> (anyrange, anymultirange)	
	&< (anyrange, anyrange)	
	&< (anyrange, anymultirange)	
	&> (anyrange, anyrange)	
	&> (anyrange, anymultirange)	
	- - (anyrange, anyrange)	
	- - (anyrange, anymultirange)	
tsquery_ops	<@ (tsquery, tsquery)	
	@> (tsquery, tsquery)	
tsvector_ops	@@ (tsvector, tsquery)	

For historical reasons, the `inet_ops` operator class is not the default class for types `inet` and `cidr`. To use it, mention the class name in `CREATE INDEX`, for example

```
CREATE INDEX ON my_table USING GIST (my_inet_column inet_ops);
```

65.2.3. Extensibility

Traditionally, implementing a new index access method meant a lot of difficult work. It was necessary to understand the inner workings of the database, such as the lock manager and Write-Ahead Log. The GiST interface has a high level of abstraction, requiring the access method implementer only to implement the semantics of the data type being accessed. The GiST layer itself takes care of concurrency, logging and searching the tree structure.

This extensibility should not be confused with the extensibility of the other standard search trees in terms of the data they can handle. For example, Postgres Pro supports extensible B-trees and hash indexes. That means that you can use Postgres Pro to build a B-tree or hash over any data type you want. But B-trees only support range predicates (`<`, `=`, `>`), and hash indexes only support equality queries.

So if you index, say, an image collection with a Postgres Pro B-tree, you can only issue queries such as “is `imagex` equal to `imagey`”, “is `imagex` less than `imagey`” and “is `imagex` greater than `imagey`”. Depending on how you define “equals”, “less than” and “greater than” in this context, this could be useful. However, by using a GiST based index, you could create ways to ask domain-specific questions, perhaps “find all images of horses” or “find all over-exposed images”.

All it takes to get a GiST access method up and running is to implement several user-defined methods, which define the behavior of keys in the tree. Of course these methods have to be pretty fancy to support fancy queries, but for all the standard queries (B-trees, R-trees, etc.) they're relatively straightforward. In short, GiST combines extensibility along with generality, code reuse, and a clean interface.

There are five methods that an index operator class for GiST must provide, and six that are optional. Correctness of the index is ensured by proper implementation of the `same`, `consistent` and `union` methods, while efficiency (size and speed) of the index will depend on the `penalty` and `picksplit` methods. Two optional methods are `compress` and `decompress`, which allow an index to have internal tree data of a different type than the data it indexes. The leaves are to be of the indexed data type, while the other tree nodes can be of any C struct (but you still have to follow Postgres Pro data type rules here, see about `varlena` for variable sized data). If the tree's internal data type exists at the SQL level, the `STORAGE` option of the `CREATE OPERATOR CLASS` command can be used. The optional eighth method is `distance`,

which is needed if the operator class wishes to support ordered scans (nearest-neighbor searches). The optional ninth method `fetch` is needed if the operator class wishes to support index-only scans, except when the `compress` method is omitted. The optional tenth method `options` is needed if the operator class has user-specified parameters. The optional eleventh method `sortsupport` is used to speed up building a GiST index.

`consistent`

Given an index entry `p` and a query value `q`, this function determines whether the index entry is “consistent” with the query; that is, could the predicate “*indexed_column indexable_operator q*” be true for any row represented by the index entry? For a leaf index entry this is equivalent to testing the indexable condition, while for an internal tree node this determines whether it is necessary to scan the subtree of the index represented by the tree node. When the result is `true`, a `recheck` flag must also be returned. This indicates whether the predicate is certainly true or only possibly true. If `recheck = false` then the index has tested the predicate condition exactly, whereas if `recheck = true` the row is only a candidate match. In that case the system will automatically evaluate the *indexable_operator* against the actual row value to see if it is really a match. This convention allows GiST to support both lossless and lossy index structures.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_consistent(internal, data_type, smallint, oid,
    internal)
RETURNS bool
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_consistent);
```

```
Datum
```

```
my_consistent(PG_FUNCTION_ARGS)
```

```
{
    GISTENTRY *entry = (GISTENTRY *) PG_GETARG_POINTER(0);
    data_type *query = PG_GETARG_DATA_TYPE_P(1);
    StrategyNumber strategy = (StrategyNumber) PG_GETARG_UINT16(2);
    /* Oid subtype = PG_GETARG_OID(3); */
    bool *recheck = (bool *) PG_GETARG_POINTER(4);
    data_type *key = DatumGetDataTypes(entry->key);
    bool retval;

    /*
     * determine return value as a function of strategy, key and query.
     *
     * Use GIST_LEAF(entry) to know where you're called in the index tree,
     * which comes handy when supporting the = operator for example (you could
     * check for non empty union() in non-leaf nodes and equality in leaf
     * nodes).
     */

    *recheck = true; /* or false if check is exact */

    PG_RETURN_BOOL(retval);
}
```

Here, `key` is an element in the index and `query` the value being looked up in the index. The `StrategyNumber` parameter indicates which operator of your operator class is being applied — it matches one of the operator numbers in the `CREATE OPERATOR CLASS` command.

Depending on which operators you have included in the class, the data type of `query` could vary with the operator, since it will be whatever type is on the right-hand side of the operator, which might

be different from the indexed data type appearing on the left-hand side. (The above code skeleton assumes that only one type is possible; if not, fetching the `query` argument value would have to depend on the operator.) It is recommended that the SQL declaration of the `consistent` function use the opclass's indexed data type for the `query` argument, even though the actual type might be something else depending on the operator.

`union`

This method consolidates information in the tree. Given a set of entries, this function generates a new index entry that represents all the given entries.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_union(internal, internal)
RETURNS storage_type
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_union);

Datum
my_union(PG_FUNCTION_ARGS)
{
    GistEntryVector *entryvec = (GistEntryVector *) PG_GETARG_POINTER(0);
    GISTENTRY *ent = entryvec->vector;
    data_type *out,
               *tmp,
               *old;
    int numranges,
        i = 0;

    numranges = entryvec->n;
    tmp = DatumGetDataType(ent[0].key);
    out = tmp;

    if (numranges == 1)
    {
        out = data_type_deep_copy(tmp);

        PG_RETURN_DATA_TYPE_P(out);
    }

    for (i = 1; i < numranges; i++)
    {
        old = out;
        tmp = DatumGetDataType(ent[i].key);
        out = my_union_implementation(out, tmp);
    }

    PG_RETURN_DATA_TYPE_P(out);
}
```

As you can see, in this skeleton we're dealing with a data type where `union(X, Y, Z) = union(union(X, Y), Z)`. It's easy enough to support data types where this is not the case, by implementing the proper union algorithm in this GiST support method.

The result of the `union` function must be a value of the index's storage type, whatever that is (it might or might not be different from the indexed column's type). The `union` function should return

a pointer to newly `palloc()`ed memory. You can't just return the input value as-is, even if there is no type change.

As shown above, the `union` function's first internal argument is actually a `GistEntryVector` pointer. The second argument is a pointer to an integer variable, which can be ignored. (It used to be required that the `union` function store the size of its result value into that variable, but this is no longer necessary.)

`compress`

Converts a data item into a format suitable for physical storage in an index page. If the `compress` method is omitted, data items are stored in the index without modification.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_compress(internal)
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_compress);
```

```
Datum
```

```
my_compress(PG_FUNCTION_ARGS)
```

```
{
    GISTENTRY *entry = (GISTENTRY *) PG_GETARG_POINTER(0);
    GISTENTRY *retval;

    if (entry->leafkey)
    {
        /* replace entry->key with a compressed version */
        compressed_data_type *compressed_data =
palloc(sizeof(compressed_data_type));

        /* fill *compressed_data from entry->key ... */

        retval = palloc(sizeof(GISTENTRY));
        gistentryinit(*retval, PointerGetDatum(compressed_data),
                      entry->rel, entry->page, entry->offset, FALSE);
    }
    else
    {
        /* typically we needn't do anything with non-leaf entries */
        retval = entry;
    }

    PG_RETURN_POINTER(retval);
}
```

You have to adapt `compressed_data_type` to the specific type you're converting to in order to compress your leaf nodes, of course.

`decompress`

Converts the stored representation of a data item into a format that can be manipulated by the other GiST methods in the operator class. If the `decompress` method is omitted, it is assumed that the other GiST methods can work directly on the stored data format. (`decompress` is not necessarily the reverse of the `compress` method; in particular, if `compress` is lossy then it's impossible for `decompress` to exactly reconstruct the original data. `decompress` is not necessarily equivalent to `fetch`, either, since the other GiST methods might not require full reconstruction of the data.)

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_decompress(internal)
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_decompress);

Datum
my_decompress(PG_FUNCTION_ARGS)
{
    PG_RETURN_POINTER(PG_GETARG_POINTER(0));
}
```

The above skeleton is suitable for the case where no decompression is needed. (But, of course, omitting the method altogether is even easier, and is recommended in such cases.)

penalty

Returns a value indicating the “cost” of inserting the new entry into a particular branch of the tree. Items will be inserted down the path of least `penalty` in the tree. Values returned by `penalty` should be non-negative. If a negative value is returned, it will be treated as zero.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_penalty(internal, internal, internal)
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT; -- in some cases penalty functions need not be strict
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_penalty);

Datum
my_penalty(PG_FUNCTION_ARGS)
{
    GISTENTRY *origentry = (GISTENTRY *) PG_GETARG_POINTER(0);
    GISTENTRY *newentry = (GISTENTRY *) PG_GETARG_POINTER(1);
    float *penalty = (float *) PG_GETARG_POINTER(2);
    data_type *orig = DatumGetDataType(origentry->key);
    data_type *new = DatumGetDataType(newentry->key);

    *penalty = my_penalty_implementation(orig, new);
    PG_RETURN_POINTER(penalty);
}
```

For historical reasons, the `penalty` function doesn't just return a `float` result; instead it has to store the value at the location indicated by the third argument. The return value per se is ignored, though it's conventional to pass back the address of that argument.

The `penalty` function is crucial to good performance of the index. It'll get used at insertion time to determine which branch to follow when choosing where to add the new entry in the tree. At query time, the more balanced the index, the quicker the lookup.

picksplit

When an index page split is necessary, this function decides which entries on the page are to stay on the old page, and which are to move to the new page.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_picksplit(internal, internal)
```

```
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_picksplit);

Datum
my_picksplit(PG_FUNCTION_ARGS)
{
    GistEntryVector *entryvec = (GistEntryVector *) PG_GETARG_POINTER(0);
    GIST_SPLITVEC *v = (GIST_SPLITVEC *) PG_GETARG_POINTER(1);
    OffsetNumber maxoff = entryvec->n - 1;
    GISTENTRY *ent = entryvec->vector;
    int i,
        nbytes;
    OffsetNumber *left,
        *right;
    data_type *tmp_union;
    data_type *unionL;
    data_type *unionR;
    GISTENTRY **raw_entryvec;

    maxoff = entryvec->n - 1;
    nbytes = (maxoff + 1) * sizeof(OffsetNumber);

    v->spl_left = (OffsetNumber *) palloc(nbytes);
    left = v->spl_left;
    v->spl_nleft = 0;

    v->spl_right = (OffsetNumber *) palloc(nbytes);
    right = v->spl_right;
    v->spl_nright = 0;

    unionL = NULL;
    unionR = NULL;

    /* Initialize the raw entry vector. */
    raw_entryvec = (GISTENTRY **) malloc(entryvec->n * sizeof(void *));
    for (i = FirstOffsetNumber; i <= maxoff; i = OffsetNumberNext(i))
        raw_entryvec[i] = &(entryvec->vector[i]);

    for (i = FirstOffsetNumber; i <= maxoff; i = OffsetNumberNext(i))
    {
        int real_index = raw_entryvec[i] - entryvec->vector;

        tmp_union = DatumGetDataType(entryvec->vector[real_index].key);
        Assert(tmp_union != NULL);

        /*
         * Choose where to put the index entries and update unionL and unionR
         * accordingly. Append the entries to either v->spl_left or
         * v->spl_right, and care about the counters.
         */

        if (my_choice_is_left(unionL, curl, unionR, curr))
        {
            if (unionL == NULL)
```

```
        unionL = tmp_union;
    else
        unionL = my_union_implementation(unionL, tmp_union);

    *left = real_index;
    ++left;
    ++(v->spl_nleft);
}
else
{
    /*
     * Same on the right
     */
}
}

v->spl_ldatum = DataTypeGetDatum(unionL);
v->spl_rdatum = DataTypeGetDatum(unionR);
PG_RETURN_POINTER(v);
}
```

Notice that the `picksplit` function's result is delivered by modifying the passed-in `v` structure. The return value per se is ignored, though it's conventional to pass back the address of `v`.

Like `penalty`, the `picksplit` function is crucial to good performance of the index. Designing suitable `penalty` and `picksplit` implementations is where the challenge of implementing well-performing GiST indexes lies.

`same`

Returns true if two index entries are identical, false otherwise. (An “index entry” is a value of the index's storage type, not necessarily the original indexed column's type.)

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_same(storage_type, storage_type, internal)
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_same);

Datum
my_same(PG_FUNCTION_ARGS)
{
    prefix_range *v1 = PG_GETARG_PREFIX_RANGE_P(0);
    prefix_range *v2 = PG_GETARG_PREFIX_RANGE_P(1);
    bool         *result = (bool *) PG_GETARG_POINTER(2);

    *result = my_eq(v1, v2);
    PG_RETURN_POINTER(result);
}
```

For historical reasons, the `same` function doesn't just return a Boolean result; instead it has to store the flag at the location indicated by the third argument. The return value per se is ignored, though it's conventional to pass back the address of that argument.

`distance`

Given an index entry `p` and a query value `q`, this function determines the index entry's “distance” from the query value. This function must be supplied if the operator class contains any ordering

operators. A query using the ordering operator will be implemented by returning index entries with the smallest “distance” values first, so the results must be consistent with the operator's semantics. For a leaf index entry the result just represents the distance to the index entry; for an internal tree node, the result must be the smallest distance that any child entry could have.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_distance(internal, data_type, smallint, oid, internal)
RETURNS float8
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

And the matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_distance);

Datum
my_distance(PG_FUNCTION_ARGS)
{
    GISTENTRY *entry = (GISTENTRY *) PG_GETARG_POINTER(0);
    data_type *query = PG_GETARG_DATA_TYPE_P(1);
    StrategyNumber strategy = (StrategyNumber) PG_GETARG_UINT16(2);
    /* Oid subtype = PG_GETARG_OID(3); */
    /* bool *recheck = (bool *) PG_GETARG_POINTER(4); */
    data_type *key = DatumGetDataTypes(entry->key);
    double      retval;

    /*
     * determine return value as a function of strategy, key and query.
     */

    PG_RETURN_FLOAT8(retval);
}
```

The arguments to the `distance` function are identical to the arguments of the `consistent` function.

Some approximation is allowed when determining the distance, so long as the result is never greater than the entry's actual distance. Thus, for example, distance to a bounding box is usually sufficient in geometric applications. For an internal tree node, the distance returned must not be greater than the distance to any of the child nodes. If the returned distance is not exact, the function must set `*recheck` to true. (This is not necessary for internal tree nodes; for them, the calculation is always assumed to be inexact.) In this case the executor will calculate the accurate distance after fetching the tuple from the heap, and reorder the tuples if necessary.

If the distance function returns `*recheck = true` for any leaf node, the original ordering operator's return type must be `float8` or `float4`, and the distance function's result values must be comparable to those of the original ordering operator, since the executor will sort using both distance function results and recalculated ordering-operator results. Otherwise, the distance function's result values can be any finite `float8` values, so long as the relative order of the result values matches the order returned by the ordering operator. (Infinity and minus infinity are used internally to handle cases such as nulls, so it is not recommended that `distance` functions return these values.)

fetch

Converts the compressed index representation of a data item into the original data type, for index-only scans. The returned data must be an exact, non-lossy copy of the originally indexed value.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_fetch(internal)
RETURNS internal
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

The argument is a pointer to a `GISTENTRY` struct. On entry, its `key` field contains a non-NULL leaf datum in compressed form. The return value is another `GISTENTRY` struct, whose `key` field contains the same datum in its original, uncompressed form. If the opclass's compress function does nothing for leaf entries, the `fetch` method can return the argument as-is. Or, if the opclass does not have a compress function, the `fetch` method can be omitted as well, since it would necessarily be a no-op.

The matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_fetch);

Datum
my_fetch(PG_FUNCTION_ARGS)
{
    GISTENTRY *entry = (GISTENTRY *) PG_GETARG_POINTER(0);
    input_data_type *in = DatumGetPointer(entry->key);
    fetched_data_type *fetched_data;
    GISTENTRY *retval;

    retval = palloc(sizeof(GISTENTRY));
    fetched_data = palloc(sizeof(fetched_data_type));

    /*
     * Convert 'fetched_data' into the a Datum of the original datatype.
     */

    /* fill *retval from fetched_data. */
    gistentryinit(*retval, PointerGetDatum(converted_datum),
                  entry->rel, entry->page, entry->offset, FALSE);

    PG_RETURN_POINTER(retval);
}
```

If the compress method is lossy for leaf entries, the operator class cannot support index-only scans, and must not define a `fetch` function.

options

Allows definition of user-visible parameters that control operator class behavior.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_options(internal)
RETURNS void
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

The function is passed a pointer to a `local_relopts` struct, which needs to be filled with a set of operator class specific options. The options can be accessed from other support functions using the `PG_HAS_OPCLASS_OPTIONS()` and `PG_GET_OPCLASS_OPTIONS()` macros.

An example implementation of `my_options()` and parameters use from other support functions are given below:

```
typedef enum MyEnumType
{
    MY_ENUM_ON,
    MY_ENUM_OFF,
    MY_ENUM_AUTO
} MyEnumType;

typedef struct
{

```

```
int32    vl_len_;    /* varlena header (do not touch directly!) */
int      int_param;  /* integer parameter */
double   real_param; /* real parameter */
MyEnumType enum_param; /* enum parameter */
int      str_param;  /* string parameter */
} MyOptionsStruct;

/* String representation of enum values */
static relopt_enum_elt_def myEnumValues[] =
{
    {"on", MY_ENUM_ON},
    {"off", MY_ENUM_OFF},
    {"auto", MY_ENUM_AUTO},
    {(const char *) NULL} /* list terminator */
};

static char *str_param_default = "default";

/*
 * Sample validator: checks that string is not longer than 8 bytes.
 */
static void
validate_my_string_relopt(const char *value)
{
    if (strlen(value) > 8)
        ereport(ERROR,
                (errcode(ERRCODE_INVALID_PARAMETER_VALUE),
                 errmsg("str_param must be at most 8 bytes")));
}

/*
 * Sample filler: switches characters to lower case.
 */
static Size
fill_my_string_relopt(const char *value, void *ptr)
{
    char    *tmp = str_tolower(value, strlen(value), DEFAULT_COLLATION_OID);
    int      len = strlen(tmp);

    if (ptr)
        strcpy((char *) ptr, tmp);

    pfree(tmp);
    return len + 1;
}

PG_FUNCTION_INFO_V1(my_options);

Datum
my_options(PG_FUNCTION_ARGS)
{
    local_relopts *relopts = (local_relopts *) PG_GETARG_POINTER(0);

    init_local_reloptions(relopts, sizeof(MyOptionsStruct));
    add_local_int_reloption(relopts, "int_param", "integer parameter",
                           100, 0, 1000000,
                           offsetof(MyOptionsStruct, int_param));
    add_local_real_reloption(relopts, "real_param", "real parameter",
```



```
        1.0, 0.0, 1000000.0,
        offsetof(MyOptionsStruct, real_param));
add_local_enum_reloption(relopts, "enum_param", "enum parameter",
        myEnumValues, MY_ENUM_ON,
        "Valid values are: \"on\", \"off\" and \"auto\".",
        offsetof(MyOptionsStruct, enum_param));
add_local_string_reloption(relopts, "str_param", "string parameter",
        str_param_default,
        &validate_my_string_relopt,
        &fill_my_string_relopt,
        offsetof(MyOptionsStruct, str_param));

PG_RETURN_VOID();
}

PG_FUNCTION_INFO_V1(my_compress);

Datum
my_compress(PG_FUNCTION_ARGS)
{
    int    int_param = 100;
    double real_param = 1.0;
    MyEnumType enum_param = MY_ENUM_ON;
    char    *str_param = str_param_default;

    /*
     * Normally, when opclass contains 'options' method, then options are always
     * passed to support functions. However, if you add 'options' method to
     * existing opclass, previously defined indexes have no options, so the
     * check is required.
     */
    if (PG_HAS_OPCLASS_OPTIONS())
    {
        MyOptionsStruct *options = (MyOptionsStruct *) PG_GET_OPCLASS_OPTIONS();

        int_param = options->int_param;
        real_param = options->real_param;
        enum_param = options->enum_param;
        str_param = GET_STRING_RELOPTION(options, str_param);
    }

    /* the rest implementation of support function */
}
```

Since the representation of the key in GiST is flexible, it may depend on user-specified parameters. For instance, the length of key signature may be specified. See `gtsvector_options()` for example.

`sortsupport`

Returns a comparator function to sort data in a way that preserves locality. It is used by `CREATE INDEX` and `REINDEX` commands. The quality of the created index depends on how well the sort order determined by the comparator function preserves locality of the inputs.

The `sortsupport` method is optional. If it is not provided, `CREATE INDEX` builds the index by inserting each tuple to the tree using the `penalty` and `picksplit` functions, which is much slower.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_sortsupport(internal)
```

```
RETURNS void
AS 'MODULE_PATHNAME'
LANGUAGE C STRICT;
```

The argument is a pointer to a `SortSupport` struct. At a minimum, the function must fill in its `comparator` field. The comparator takes three arguments: two `Datums` to compare, and a pointer to the `SortSupport` struct. The `Datums` are the two indexed values in the format that they are stored in the index; that is, in the format returned by the `compress` method. The full API is defined in [src/include/utils/sortsupport.h](#).

The matching code in the C module could then follow this skeleton:

```
PG_FUNCTION_INFO_V1(my_sortsupport);

static int
my_fastcmp(Datum x, Datum y, SortSupport ssup)
{
    /* establish order between x and y by computing some sorting value z */

    int z1 = ComputeSpatialCode(x);
    int z2 = ComputeSpatialCode(y);

    return z1 == z2 ? 0 : z1 > z2 ? 1 : -1;
}

Datum
my_sortsupport(PG_FUNCTION_ARGS)
{
    SortSupport ssup = (SortSupport) PG_GETARG_POINTER(0);

    ssup->comparator = my_fastcmp;
    PG_RETURN_VOID();
}
```

All the GiST support methods are normally called in short-lived memory contexts; that is, `CurrentMemoryContext` will get reset after each tuple is processed. It is therefore not very important to worry about `pfree`'ing everything you `palloc`. However, in some cases it's useful for a support method to cache data across repeated calls. To do that, allocate the longer-lived data in `fcinfo->flinfo->fn_mcxt`, and keep a pointer to it in `fcinfo->flinfo->fn_extra`. Such data will survive for the life of the index operation (e.g., a single GiST index scan, index build, or index tuple insertion). Be careful to `pfree` the previous value when replacing a `fn_extra` value, or the leak will accumulate for the duration of the operation.

65.2.4. Implementation

65.2.4.1. GiST Index Build Methods

The simplest way to build a GiST index is just to insert all the entries, one by one. This tends to be slow for large indexes, because if the index tuples are scattered across the index and the index is large enough to not fit in cache, a lot of random I/O will be needed. Postgres Pro supports two alternative methods for initial build of a GiST index: *sorted* and *buffered* modes.

The sorted method is only available if each of the opclasses used by the index provides a `sortsupport` function, as described in [Section 65.2.3](#). If they do, this method is usually the best, so it is used by default.

The buffered method works by not inserting tuples directly into the index right away. It can dramatically reduce the amount of random I/O needed for non-ordered data sets. For well-ordered data sets the benefit is smaller or non-existent, because only a small number of pages receive new tuples at a time, and those pages fit in cache even if the index as a whole does not.

The buffered method needs to call the `penalty` function more often than the simple method does, which consumes some extra CPU resources. Also, the buffers need temporary disk space, up to the size of

the resulting index. Buffering can also influence the quality of the resulting index, in both positive and negative directions. That influence depends on various factors, like the distribution of the input data and the operator class implementation.

If sorting is not possible, then by default a GiST index build switches to the buffering method when the index size reaches [effective_cache_size](#). Buffering can be manually forced or prevented by the `buffering` parameter to the `CREATE INDEX` command. The default behavior is good for most cases, but turning buffering off might speed up the build somewhat if the input data is ordered.

65.2.5. Examples

The Postgres Pro core system currently provides text search support (indexing for `tsvector` and `tsquery`) as well as R-Tree equivalent functionality for some of the built-in geometric data types. The following `contrib` modules also contain GiST operator classes:

`btree_gist`

B-tree equivalent functionality for several data types

`cube`

Indexing for multidimensional cubes

`hstore`

Module for storing (key, value) pairs

`intarray`

RD-Tree for one-dimensional array of `int4` values

`ltree`

Indexing for tree-like structures

`pg_trgm`

Text similarity using trigram matching

`seg`

Indexing for “float ranges”

65.3. SP-GiST Indexes

65.3.1. Introduction

SP-GiST is an abbreviation for space-partitioned GiST. SP-GiST supports partitioned search trees, which facilitate development of a wide range of different non-balanced data structures, such as quad-trees, k-d trees, and radix trees (tries). The common feature of these structures is that they repeatedly divide the search space into partitions that need not be of equal size. Searches that are well matched to the partitioning rule can be very fast.

These popular data structures were originally developed for in-memory usage. In main memory, they are usually designed as a set of dynamically allocated nodes linked by pointers. This is not suitable for direct storing on disk, since these chains of pointers can be rather long which would require too many disk accesses. In contrast, disk-based data structures should have a high fanout to minimize I/O. The challenge addressed by SP-GiST is to map search tree nodes to disk pages in such a way that a search need access only a few disk pages, even if it traverses many nodes.

Like GiST, SP-GiST is meant to allow the development of custom data types with the appropriate access methods, by an expert in the domain of the data type, rather than a database expert.

Some of the information here is derived from Purdue University's SP-GiST Indexing Project [web site](#). The SP-GiST implementation in Postgres Pro is primarily maintained by Teodor Sigaev and Oleg Bartunov, and there is more information on their [web site](#).

65.3.2. Built-in Operator Classes

The core Postgres Pro distribution includes the SP-GiST operator classes shown in [Table 65.2](#).

Table 65.2. Built-in SP-GiST Operator Classes

Name	Indexable Operators	Ordering Operators
box_ops	<< (box,box)	<-> (box,point)
	&< (box,box)	
	&> (box,box)	
	>> (box,box)	
	<@ (box,box)	
	@> (box,box)	
	~= (box,box)	
	&& (box,box)	
	<< (box,box)	
	&< (box,box)	
	&> (box,box)	
	>> (box,box)	
inet_ops	<< (inet,inet)	
	<<= (inet,inet)	
	>> (inet,inet)	
	>>= (inet,inet)	
	= (inet,inet)	
	<> (inet,inet)	
	< (inet,inet)	
	<= (inet,inet)	
	> (inet,inet)	
	>= (inet,inet)	
	&& (inet,inet)	
kd_point_ops	>> (point,point)	<-> (point,point)
	<< (point,point)	
	>> (point,point)	
	<< (point,point)	
	~= (point,point)	
	<@ (point,box)	
poly_ops	<< (polygon,polygon)	<-> (polygon,point)
	&< (polygon,polygon)	
	&> (polygon,polygon)	
	>> (polygon,polygon)	
	<@ (polygon,polygon)	
	@> (polygon,polygon)	
	~= (polygon,polygon)	

Name	Indexable Operators	Ordering Operators
	&& (polygon,polygon)	
	<< (polygon,polygon)	
	&< (polygon,polygon)	
	>> (polygon,polygon)	
	&> (polygon,polygon)	
quad_point_ops	>> (point,point)	<-> (point,point)
	<< (point,point)	
	>> (point,point)	
	<< (point,point)	
	~= (point,point)	
	<@ (point,box)	
range_ops	= (anyrange,anyrange)	
	&& (anyrange,anyrange)	
	@> (anyrange,anyelement)	
	@> (anyrange,anyrange)	
	<@ (anyrange,anyrange)	
	<< (anyrange,anyrange)	
	>> (anyrange,anyrange)	
	&< (anyrange,anyrange)	
	&> (anyrange,anyrange)	
	- - (anyrange,anyrange)	
text_ops	= (text,text)	
	< (text,text)	
	<= (text,text)	
	> (text,text)	
	>= (text,text)	
	~<~ (text,text)	
	~<=~ (text,text)	
	~>=~ (text,text)	
	~>~ (text,text)	
	^@ (text,text)	

Of the two operator classes for type `point`, `quad_point_ops` is the default. `kd_point_ops` supports the same operators but uses a different index data structure that may offer better performance in some applications.

The `quad_point_ops`, `kd_point_ops` and `poly_ops` operator classes support the `<->` ordering operator, which enables the k-nearest neighbor (k-NN) search over indexed point or polygon data sets.

65.3.3. Extensibility

SP-GiST offers an interface with a high level of abstraction, requiring the access method developer to implement only methods specific to a given data type. The SP-GiST core is responsible for efficient disk mapping and searching the tree structure. It also takes care of concurrency and logging considerations.

Leaf tuples of an SP-GiST tree usually contain values of the same data type as the indexed column, although it is also possible for them to contain lossy representations of the indexed column. Leaf tuples stored at the root level will directly represent the original indexed data value, but leaf tuples at lower levels might contain only a partial value, such as a suffix. In that case the operator class support functions must be able to reconstruct the original value using information accumulated from the inner tuples that are passed through to reach the leaf level.

When an SP-GiST index is created with `INCLUDE` columns, the values of those columns are also stored in leaf tuples. The `INCLUDE` columns are of no concern to the SP-GiST operator class, so they are not discussed further here.

Inner tuples are more complex, since they are branching points in the search tree. Each inner tuple contains a set of one or more *nodes*, which represent groups of similar leaf values. A node contains a downlink that leads either to another, lower-level inner tuple, or to a short list of leaf tuples that all lie on the same index page. Each node normally has a *label* that describes it; for example, in a radix tree the node label could be the next character of the string value. (Alternatively, an operator class can omit the node labels, if it works with a fixed set of nodes for all inner tuples; see [Section 65.3.4.2](#).) Optionally, an inner tuple can have a *prefix* value that describes all its members. In a radix tree this could be the common prefix of the represented strings. The prefix value is not necessarily really a prefix, but can be any data needed by the operator class; for example, in a quad-tree it can store the central point that the four quadrants are measured with respect to. A quad-tree inner tuple would then also contain four nodes corresponding to the quadrants around this central point.

Some tree algorithms require knowledge of level (or depth) of the current tuple, so the SP-GiST core provides the possibility for operator classes to manage level counting while descending the tree. There is also support for incrementally reconstructing the represented value when that is needed, and for passing down additional data (called *traverse values*) during a tree descent.

Note

The SP-GiST core code takes care of null entries. Although SP-GiST indexes do store entries for nulls in indexed columns, this is hidden from the index operator class code: no null index entries or search conditions will ever be passed to the operator class methods. (It is assumed that SP-GiST operators are strict and so cannot succeed for null values.) Null values are therefore not discussed further here.

There are five user-defined methods that an index operator class for SP-GiST must provide, and two are optional. All five mandatory methods follow the convention of accepting two `internal` arguments, the first of which is a pointer to a C struct containing input values for the support method, while the second argument is a pointer to a C struct where output values must be placed. Four of the mandatory methods just return `void`, since all their results appear in the output struct; but `leaf_consistent` returns a `boolean` result. The methods must not modify any fields of their input structs. In all cases, the output struct is initialized to zeroes before calling the user-defined method. The optional sixth method `compress` accepts a `datum` to be indexed as the only argument and returns a value suitable for physical storage in a leaf tuple. The optional seventh method `options` accepts an `internal` pointer to a C struct, where opclass-specific parameters should be placed, and returns `void`.

The five mandatory user-defined methods are:

`config`

Returns static information about the index implementation, including the data type OIDs of the prefix and node label data types.

The SQL declaration of the function must look like this:

```
CREATE FUNCTION my_config(internal, internal) RETURNS void ...
```

The first argument is a pointer to a `spgConfigIn` C struct, containing input data for the function. The second argument is a pointer to a `spgConfigOut` C struct, which the function must fill with result data.

```
typedef struct spgConfigIn
{
    Oid          attType;          /* Data type to be indexed */
} spgConfigIn;

typedef struct spgConfigOut
{
    Oid          prefixType;       /* Data type of inner-tuple prefixes */
    Oid          labelType;       /* Data type of inner-tuple node labels */
    Oid          leafType;        /* Data type of leaf-tuple values */
    bool         canReturnData;   /* Opclass can reconstruct original data */
    bool         longValuesOK;    /* Opclass can cope with values > 1 page */
} spgConfigOut;
```

`attType` is passed in order to support polymorphic index operator classes; for ordinary fixed-data-type operator classes, it will always have the same value and so can be ignored.

For operator classes that do not use prefixes, `prefixType` can be set to `VOIDOID`. Likewise, for operator classes that do not use node labels, `labelType` can be set to `VOIDOID`. `canReturnData` should be set true if the operator class is capable of reconstructing the originally-supplied index value. `longValuesOK` should be set true only when the `attType` is of variable length and the operator class is capable of segmenting long values by repeated suffixing (see [Section 65.3.4.1](#)).

`leafType` should match the index storage type defined by the operator class's `opkeytype` catalog entry. (Note that `opkeytype` can be zero, implying the storage type is the same as the operator class's input type, which is the most common situation.) For reasons of backward compatibility, the `config` method can set `leafType` to some other value, and that value will be used; but this is deprecated since the index contents are then incorrectly identified in the catalogs. Also, it's permissible to leave `leafType` uninitialized (zero); that is interpreted as meaning the index storage type derived from `opkeytype`.

When `attType` and `leafType` are different, the optional method `compress` must be provided. Method `compress` is responsible for transformation of datums to be indexed from `attType` to `leafType`.

`choose`

Chooses a method for inserting a new value into an inner tuple.

The SQL declaration of the function must look like this:

```
CREATE FUNCTION my_choose(internal, internal) RETURNS void ...
```

The first argument is a pointer to a `spgChooseIn` C struct, containing input data for the function. The second argument is a pointer to a `spgChooseOut` C struct, which the function must fill with result data.

```
typedef struct spgChooseIn
{
    Datum        datum;           /* original datum to be indexed */
    Datum        leafDatum;       /* current datum to be stored at leaf */
    int          level;           /* current level (counting from zero) */

    /* Data from current inner tuple */
    bool         allTheSame;      /* tuple is marked all-the-same? */
    bool         hasPrefix;       /* tuple has a prefix? */
    Datum        prefixDatum;     /* if so, the prefix value */
    int          nNodes;          /* number of nodes in the inner tuple */
    Datum        *nodeLabels;     /* node label values (NULL if none) */
} spgChooseIn;
```

```
typedef enum spgChooseResultType
{
    spgMatchNode = 1,          /* descend into existing node */
    spgAddNode,                /* add a node to the inner tuple */
    spgSplitTuple              /* split inner tuple (change its prefix) */
} spgChooseResultType;

typedef struct spgChooseOut
{
    spgChooseResultType resultType;    /* action code, see above */
    union
    {
        struct                /* results for spgMatchNode */
        {
            int                nodeN;    /* descend to this node (index from 0) */
            int                levelAdd; /* increment level by this much */
            Datum              restDatum; /* new leaf datum */
        } matchNode;
        struct                /* results for spgAddNode */
        {
            Datum              nodeLabel; /* new node's label */
            int                nodeN;    /* where to insert it (index from 0) */
        } addNode;
        struct                /* results for spgSplitTuple */
        {
            /* Info to form new upper-level inner tuple with one child tuple */
            bool                prefixHasPrefix; /* tuple should have a prefix? */
            Datum              prefixPrefixDatum; /* if so, its value */
            int                prefixNNodes;    /* number of nodes */
            Datum              *prefixNodeLabels; /* their labels (or NULL for
                                                    * no labels) */
            int                childNodeN;    /* which node gets child tuple */

            /* Info to form new lower-level inner tuple with all old nodes */
            bool                postfixHasPrefix; /* tuple should have a prefix? */
            Datum              postfixPrefixDatum; /* if so, its value */
        } splitTuple;
    } result;
} spgChooseOut;
```

`datum` is the original datum of `spgConfigIn.attType` type that was to be inserted into the index. `leafDatum` is a value of `spgConfigOut.leafType` type, which is initially a result of method `compress` applied to `datum` when method `compress` is provided, or the same value as `datum` otherwise. `leafDatum` can change at lower levels of the tree if the `choose` or `picksplit` methods change it. When the insertion search reaches a leaf page, the current value of `leafDatum` is what will be stored in the newly created leaf tuple. `level` is the current inner tuple's level, starting at zero for the root level. `allTheSame` is true if the current inner tuple is marked as containing multiple equivalent nodes (see [Section 65.3.4.3](#)). `hasPrefix` is true if the current inner tuple contains a prefix; if so, `prefixDatum` is its value. `nNodes` is the number of child nodes contained in the inner tuple, and `nodeLabels` is an array of their label values, or NULL if there are no labels.

The `choose` function can determine either that the new value matches one of the existing child nodes, or that a new child node must be added, or that the new value is inconsistent with the tuple prefix and so the inner tuple must be split to create a less restrictive prefix.

If the new value matches one of the existing child nodes, set `resultType` to `spgMatchNode`. Set `nodeN` to the index (from zero) of that node in the node array. Set `levelAdd` to the increment in `level` caused by descending through that node, or leave it as zero if the operator class does not use levels.

Set `restDatum` to equal `leafDatum` if the operator class does not modify datums from one level to the next, or otherwise set it to the modified value to be used as `leafDatum` at the next level.

If a new child node must be added, set `resultType` to `spgAddNode`. Set `nodeLabel` to the label to be used for the new node, and set `nodeN` to the index (from zero) at which to insert the node in the node array. After the node has been added, the `choose` function will be called again with the modified inner tuple; that call should result in an `spgMatchNode` result.

If the new value is inconsistent with the tuple prefix, set `resultType` to `spgSplitTuple`. This action moves all the existing nodes into a new lower-level inner tuple, and replaces the existing inner tuple with a tuple having a single downlink pointing to the new lower-level inner tuple. Set `prefixHasPrefix` to indicate whether the new upper tuple should have a prefix, and if so set `prefixPrefixDatum` to the prefix value. This new prefix value must be sufficiently less restrictive than the original to accept the new value to be indexed. Set `prefixNNodes` to the number of nodes needed in the new tuple, and set `prefixNodeLabels` to a palloc'd array holding their labels, or to NULL if node labels are not required. Note that the total size of the new upper tuple must be no more than the total size of the tuple it is replacing; this constrains the lengths of the new prefix and new labels. Set `childNodeN` to the index (from zero) of the node that will downlink to the new lower-level inner tuple. Set `postfixHasPrefix` to indicate whether the new lower-level inner tuple should have a prefix, and if so set `postfixPrefixDatum` to the prefix value. The combination of these two prefixes and the downlink node's label (if any) must have the same meaning as the original prefix, because there is no opportunity to alter the node labels that are moved to the new lower-level tuple, nor to change any child index entries. After the node has been split, the `choose` function will be called again with the replacement inner tuple. That call may return an `spgAddNode` result, if no suitable node was created by the `spgSplitTuple` action. Eventually `choose` must return `spgMatchNode` to allow the insertion to descend to the next level.

`picksplit`

Decides how to create a new inner tuple over a set of leaf tuples.

The SQL declaration of the function must look like this:

```
CREATE FUNCTION my_picksplit(internal, internal) RETURNS void ...
```

The first argument is a pointer to a `spgPickSplitIn` C struct, containing input data for the function. The second argument is a pointer to a `spgPickSplitOut` C struct, which the function must fill with result data.

```
typedef struct spgPickSplitIn
{
    int          nTuples;          /* number of leaf tuples */
    Datum        *datums;          /* their datums (array of length nTuples) */
    int          level;           /* current level (counting from zero) */
} spgPickSplitIn;

typedef struct spgPickSplitOut
{
    bool          hasPrefix;       /* new inner tuple should have a prefix? */
    Datum        prefixDatum;     /* if so, its value */

    int          nNodes;          /* number of nodes for new inner tuple */
    Datum        *nodeLabels;     /* their labels (or NULL for no labels) */

    int          *mapTuplesToNodes; /* node index for each leaf tuple */
    Datum        *leafTupleDatums; /* datum to store in each new leaf tuple */
} spgPickSplitOut;
```

`nTuples` is the number of leaf tuples provided. `datums` is an array of their datum values of `spgConfigOut.leafType` type. `level` is the current level that all the leaf tuples share, which will become the level of the new inner tuple.

Set `hasPrefix` to indicate whether the new inner tuple should have a prefix, and if so set `prefixDatum` to the prefix value. Set `nNodes` to indicate the number of nodes that the new inner tuple will contain, and set `nodeLabels` to an array of their label values, or to NULL if node labels are not required. Set `mapTuplesToNodes` to an array that gives the index (from zero) of the node that each leaf tuple should be assigned to. Set `leafTupleDatums` to an array of the values to be stored in the new leaf tuples (these will be the same as the input `datums` if the operator class does not modify datums from one level to the next). Note that the `picksplit` function is responsible for palloc'ing the `nodeLabels`, `mapTuplesToNodes` and `leafTupleDatums` arrays.

If more than one leaf tuple is supplied, it is expected that the `picksplit` function will classify them into more than one node; otherwise it is not possible to split the leaf tuples across multiple pages, which is the ultimate purpose of this operation. Therefore, if the `picksplit` function ends up placing all the leaf tuples in the same node, the core SP-GiST code will override that decision and generate an inner tuple in which the leaf tuples are assigned at random to several identically-labeled nodes. Such a tuple is marked `allTheSame` to signify that this has happened. The `choose` and `inner_consistent` functions must take suitable care with such inner tuples. See [Section 65.3.4.3](#) for more information.

`picksplit` can be applied to a single leaf tuple only in the case that the `config` function set `long-ValuesOK` to true and a larger-than-a-page input value has been supplied. In this case the point of the operation is to strip off a prefix and produce a new, shorter leaf datum value. The call will be repeated until a leaf datum short enough to fit on a page has been produced. See [Section 65.3.4.1](#) for more information.

`inner_consistent`

Returns set of nodes (branches) to follow during tree search.

The SQL declaration of the function must look like this:

```
CREATE FUNCTION my_inner_consistent(internal, internal) RETURNS void ...
```

The first argument is a pointer to a `spgInnerConsistentIn` C struct, containing input data for the function. The second argument is a pointer to a `spgInnerConsistentOut` C struct, which the function must fill with result data.

```
typedef struct spgInnerConsistentIn
{
    ScanKey      scankeys;          /* array of operators and comparison values */
    ScanKey      orderbys;          /* array of ordering operators and comparison
                                     * values */
    int          nkeys;             /* length of scankeys array */
    int          norderbys;         /* length of orderbys array */

    Datum        reconstructedValue; /* value reconstructed at parent */
    void         *traversalValue;    /* opclass-specific traverse value */
    MemoryContext traversalMemoryContext; /* put new traverse values here */
    int          level;             /* current level (counting from zero) */
    bool         returnData;        /* original data must be returned? */

    /* Data from current inner tuple */
    bool         allTheSame;        /* tuple is marked all-the-same? */
    bool         hasPrefix;         /* tuple has a prefix? */
    Datum        prefixDatum;       /* if so, the prefix value */
    int          nNodes;            /* number of nodes in the inner tuple */
    Datum        *nodeLabels;       /* node label values (NULL if none) */
} spgInnerConsistentIn;

typedef struct spgInnerConsistentOut
{
    int          nNodes;            /* number of child nodes to be visited */
    int          *nodeNumbers;      /* their indexes in the node array */
}
```

```

int          *levelAdds;          /* increment level by this much for each */
Datum        *reconstructedValues; /* associated reconstructed values */
void         **traversalValues;   /* opclass-specific traverse values */
double       **distances;         /* associated distances */
} spgInnerConsistentOut;

```

The array `scankeys`, of length `nkeys`, describes the index search condition(s). These conditions are combined with AND — only index entries that satisfy all of them are interesting. (Note that `nkeys = 0` implies that all index entries satisfy the query.) Usually the consistent function only cares about the `sk_strategy` and `sk_argument` fields of each array entry, which respectively give the indexable operator and comparison value. In particular it is not necessary to check `sk_flags` to see if the comparison value is NULL, because the SP-GiST core code will filter out such conditions. The array `orderbys`, of length `norderbys`, describes ordering operators (if any) in the same manner. `reconstructedValue` is the value reconstructed for the parent tuple; it is (Datum) 0 at the root level or if the `inner_consistent` function did not provide a value at the parent level. `traversalValue` is a pointer to any traverse data passed down from the previous call of `inner_consistent` on the parent index tuple, or NULL at the root level. `traversalMemoryContext` is the memory context in which to store output traverse values (see below). `level` is the current inner tuple's level, starting at zero for the root level. `returnData` is true if reconstructed data is required for this query; this will only be so if the `config` function asserted `canReturnData`. `allTheSame` is true if the current inner tuple is marked “all-the-same”; in this case all the nodes have the same label (if any) and so either all or none of them match the query (see [Section 65.3.4.3](#)). `hasPrefix` is true if the current inner tuple contains a prefix; if so, `prefixDatum` is its value. `nNodes` is the number of child nodes contained in the inner tuple, and `nodeLabels` is an array of their label values, or NULL if the nodes do not have labels.

`nNodes` must be set to the number of child nodes that need to be visited by the search, and `nodeNumbers` must be set to an array of their indexes. If the operator class keeps track of levels, set `levelAdds` to an array of the level increments required when descending to each node to be visited. (Often these increments will be the same for all the nodes, but that's not necessarily so, so an array is used.) If value reconstruction is needed, set `reconstructedValues` to an array of the values reconstructed for each child node to be visited; otherwise, leave `reconstructedValues` as NULL. The reconstructed values are assumed to be of type `spgConfigOut.leafType`. (However, since the core system will do nothing with them except possibly copy them, it is sufficient for them to have the same `typlen` and `typbyval` properties as `leafType`.) If ordered search is performed, set `distances` to an array of distance values according to `orderbys` array (nodes with lowest distances will be processed first). Leave it NULL otherwise. If it is desired to pass down additional out-of-band information (“traverse values”) to lower levels of the tree search, set `traversalValues` to an array of the appropriate traverse values, one for each child node to be visited; otherwise, leave `traversalValues` as NULL. Note that the `inner_consistent` function is responsible for `palloc`'ing the `nodeNumbers`, `levelAdds`, `distances`, `reconstructedValues`, and `traversalValues` arrays in the current memory context. However, any output traverse values pointed to by the `traversalValues` array should be allocated in `traversalMemoryContext`. Each traverse value must be a single `palloc`'d chunk.

`leaf_consistent`

Returns true if a leaf tuple satisfies a query.

The SQL declaration of the function must look like this:

```
CREATE FUNCTION my_leaf_consistent(internal, internal) RETURNS bool ...
```

The first argument is a pointer to a `spgLeafConsistentIn` C struct, containing input data for the function. The second argument is a pointer to a `spgLeafConsistentOut` C struct, which the function must fill with result data.

```

typedef struct spgLeafConsistentIn
{
    ScanKey      scankeys;          /* array of operators and comparison values */
    ScanKey      orderbys;          /* array of ordering operators and comparison
                                   * values */
    int          nkeys;             /* length of scankeys array */

```

```
int            norderbys;          /* length of orderbys array */

Datum          reconstructedValue; /* value reconstructed at parent */
void          *traversalValue; /* opclass-specific traverse value */
int           level;             /* current level (counting from zero) */
bool          returnData;        /* original data must be returned? */

Datum          leafDatum;          /* datum in leaf tuple */
} spgLeafConsistentIn;

typedef struct spgLeafConsistentOut
{
    Datum          leafValue;        /* reconstructed original data, if any */
    bool          recheck;           /* set true if operator must be rechecked */
    bool          recheckDistances; /* set true if distances must be rechecked */
    double        *distances;        /* associated distances */
} spgLeafConsistentOut;
```

The array `scankeys`, of length `nkeys`, describes the index search condition(s). These conditions are combined with AND — only index entries that satisfy all of them satisfy the query. (Note that `nkeys = 0` implies that all index entries satisfy the query.) Usually the consistent function only cares about the `sk_strategy` and `sk_argument` fields of each array entry, which respectively give the indexable operator and comparison value. In particular it is not necessary to check `sk_flags` to see if the comparison value is NULL, because the SP-GiST core code will filter out such conditions. The array `orderbys`, of length `norderbys`, describes the ordering operators in the same manner. `reconstructedValue` is the value reconstructed for the parent tuple; it is (Datum) 0 at the root level or if the `inner_consistent` function did not provide a value at the parent level. `traversalValue` is a pointer to any traverse data passed down from the previous call of `inner_consistent` on the parent index tuple, or NULL at the root level. `level` is the current leaf tuple's level, starting at zero for the root level. `returnData` is true if reconstructed data is required for this query; this will only be so if the `config` function asserted `canReturnData`. `leafDatum` is the key value of `spgConfigOut.leafType` stored in the current leaf tuple.

The function must return `true` if the leaf tuple matches the query, or `false` if not. In the `true` case, if `returnData` is true then `leafValue` must be set to the value (of type `spgConfigIn.attType`) originally supplied to be indexed for this leaf tuple. Also, `recheck` may be set to `true` if the match is uncertain and so the operator(s) must be re-applied to the actual heap tuple to verify the match. If ordered search is performed, set `distances` to an array of distance values according to `orderbys` array. Leave it NULL otherwise. If at least one of returned distances is not exact, set `recheckDistances` to `true`. In this case, the executor will calculate the exact distances after fetching the tuple from the heap, and will reorder the tuples if needed.

The optional user-defined methods are:

`Datum compress(Datum in)`

Converts a data item into a format suitable for physical storage in a leaf tuple of the index. It accepts a value of type `spgConfigIn.attType` and returns a value of type `spgConfigOut.leafType`. The output value must not contain an out-of-line TOAST pointer.

Note: the `compress` method is only applied to values to be stored. The consistent methods receive query `scankeys` unchanged, without transformation using `compress`.

`options`

Defines a set of user-visible parameters that control operator class behavior.

The SQL declaration of the function must look like this:

```
CREATE OR REPLACE FUNCTION my_options(internal)
RETURNS void
```

```
AS 'MODULE_PATHNAME'  
LANGUAGE C STRICT;
```

The function is passed a pointer to a `local_relopts` struct, which needs to be filled with a set of operator class specific options. The options can be accessed from other support functions using the `PG_HAS_OPCLASS_OPTIONS()` and `PG_GET_OPCLASS_OPTIONS()` macros.

Since the representation of the key in SP-GiST is flexible, it may depend on user-specified parameters.

All the SP-GiST support methods are normally called in a short-lived memory context; that is, `CurrentMemoryContext` will be reset after processing of each tuple. It is therefore not very important to worry about `pfree`'ing everything you `palloc`. (The `config` method is an exception: it should try to avoid leaking memory. But usually the `config` method need do nothing but assign constants into the passed parameter struct.)

If the indexed column is of a collatable data type, the index collation will be passed to all the support methods, using the standard `PG_GET_COLLATION()` mechanism.

65.3.4. Implementation

This section covers implementation details and other tricks that are useful for implementers of SP-GiST operator classes to know.

65.3.4.1. SP-GiST Limits

Individual leaf tuples and inner tuples must fit on a single index page (8kB by default). Therefore, when indexing values of variable-length data types, long values can only be supported by methods such as radix trees, in which each level of the tree includes a prefix that is short enough to fit on a page, and the final leaf level includes a suffix also short enough to fit on a page. The operator class should set `longValuesOK` to true only if it is prepared to arrange for this to happen. Otherwise, the SP-GiST core will reject any request to index a value that is too large to fit on an index page.

Likewise, it is the operator class's responsibility that inner tuples do not grow too large to fit on an index page; this limits the number of child nodes that can be used in one inner tuple, as well as the maximum size of a prefix value.

Another limitation is that when an inner tuple's node points to a set of leaf tuples, those tuples must all be in the same index page. (This is a design decision to reduce seeking and save space in the links that chain such tuples together.) If the set of leaf tuples grows too large for a page, a split is performed and an intermediate inner tuple is inserted. For this to fix the problem, the new inner tuple *must* divide the set of leaf values into more than one node group. If the operator class's `picksplit` function fails to do that, the SP-GiST core resorts to extraordinary measures described in [Section 65.3.4.3](#).

When `longValuesOK` is true, it is expected that successive levels of the SP-GiST tree will absorb more and more information into the prefixes and node labels of the inner tuples, making the required leaf datum smaller and smaller, so that eventually it will fit on a page. To prevent bugs in operator classes from causing infinite insertion loops, the SP-GiST core will raise an error if the leaf datum does not become any smaller within ten cycles of `choose` method calls.

65.3.4.2. SP-GiST Without Node Labels

Some tree algorithms use a fixed set of nodes for each inner tuple; for example, in a quad-tree there are always exactly four nodes corresponding to the four quadrants around the inner tuple's centroid point. In such a case the code typically works with the nodes by number, and there is no need for explicit node labels. To suppress node labels (and thereby save some space), the `picksplit` function can return NULL for the `nodeLabels` array, and likewise the `choose` function can return NULL for the `prefixNodeLabels` array during a `spgSplitTuple` action. This will in turn result in `nodeLabels` being NULL during subsequent calls to `choose` and `inner_consistent`. In principle, node labels could be used for some inner tuples and omitted for others in the same index.

When working with an inner tuple having unlabeled nodes, it is an error for `choose` to return `spgAddNode`, since the set of nodes is supposed to be fixed in such cases.

65.3.4.3. “All-the-Same” Inner Tuples

The SP-GiST core can override the results of the operator class's `picksplit` function when `picksplit` fails to divide the supplied leaf values into at least two node categories. When this happens, the new inner tuple is created with multiple nodes that each have the same label (if any) that `picksplit` gave to the one node it did use, and the leaf values are divided at random among these equivalent nodes. The `allTheSame` flag is set on the inner tuple to warn the `choose` and `inner_consistent` functions that the tuple does not have the node set that they might otherwise expect.

When dealing with an `allTheSame` tuple, a `choose` result of `spgMatchNode` is interpreted to mean that the new value can be assigned to any of the equivalent nodes; the core code will ignore the supplied `nodeN` value and descend into one of the nodes at random (so as to keep the tree balanced). It is an error for `choose` to return `spgAddNode`, since that would make the nodes not all equivalent; the `spgSplitTuple` action must be used if the value to be inserted doesn't match the existing nodes.

When dealing with an `allTheSame` tuple, the `inner_consistent` function should return either all or none of the nodes as targets for continuing the index search, since they are all equivalent. This may or may not require any special-case code, depending on how much the `inner_consistent` function normally assumes about the meaning of the nodes.

65.4. GIN Indexes

65.4.1. Introduction

GIN stands for Generalized Inverted Index. GIN is designed for handling cases where the items to be indexed are composite values, and the queries to be handled by the index need to search for element values that appear within the composite items. For example, the items could be documents, and the queries could be searches for documents containing specific words.

We use the word *item* to refer to a composite value that is to be indexed, and the word *key* to refer to an element value. GIN always stores and searches for keys, not item values per se.

A GIN index stores a set of (key, posting list) pairs, where a *posting list* is a set of row IDs in which the key occurs. The same row ID can appear in multiple posting lists, since an item can contain more than one key. Each key value is stored only once, so a GIN index is very compact for cases where the same key appears many times.

GIN is generalized in the sense that the GIN access method code does not need to know the specific operations that it accelerates. Instead, it uses custom strategies defined for particular data types. The strategy defines how keys are extracted from indexed items and query conditions, and how to determine whether a row that contains some of the key values in a query actually satisfies the query.

One advantage of GIN is that it allows the development of custom data types with the appropriate access methods, by an expert in the domain of the data type, rather than a database expert. This is much the same advantage as using GiST.

The GIN implementation in Postgres Pro is primarily maintained by Teodor Sigaev and Oleg Bartunov. There is more information about GIN on their [website](#).

65.4.2. Built-in Operator Classes

The core Postgres Pro distribution includes the GIN operator classes shown in [Table 65.3](#). (Some of the optional modules described in [Appendix F](#) provide additional GIN operator classes.)

Table 65.3. Built-in GIN Operator Classes

Name	Indexable Operators
array_ops	&& (anyarray, anyarray)
	@> (anyarray, anyarray)
	<@ (anyarray, anyarray)

Name	Indexable Operators
	= (anyarray, anyarray)
jsonb_ops	@> (jsonb, jsonb)
	@? (jsonb, jsonpath)
	@@ (jsonb, jsonpath)
	? (jsonb, text)
	? (jsonb, text[])
	?& (jsonb, text[])
jsonb_path_ops	@> (jsonb, jsonb)
	@? (jsonb, jsonpath)
	@@ (jsonb, jsonpath)
tsvector_ops	@@ (tsvector, tsquery)

Of the two operator classes for type `jsonb`, `jsonb_ops` is the default. `jsonb_path_ops` supports fewer operators but offers better performance for those operators. See [Section 8.14.4](#) for details.

65.4.3. Extensibility

The GIN interface has a high level of abstraction, requiring the access method implementer only to implement the semantics of the data type being accessed. The GIN layer itself takes care of concurrency, logging and searching the tree structure.

All it takes to get a GIN access method working is to implement a few user-defined methods, which define the behavior of keys in the tree and the relationships between keys, indexed items, and indexable queries. In short, GIN combines extensibility with generality, code reuse, and a clean interface.

There are two methods that an operator class for GIN must provide:

```
Datum *extractValue(Datum itemValue, int32 *nkeys, bool **nullFlags)
```

Returns a palloc'd array of keys given an item to be indexed. The number of returned keys must be stored into `*nkeys`. If any of the keys can be null, also palloc an array of `*nkeys` `bool` fields, store its address at `*nullFlags`, and set these null flags as needed. `*nullFlags` can be left `NULL` (its initial value) if all keys are non-null. The return value can be `NULL` if the item contains no keys.

```
Datum *extractQuery(Datum query, int32 *nkeys, StrategyNumber n, bool **pmatch, Pointer
**extra_data, bool **nullFlags, int32 *searchMode)
```

Returns a palloc'd array of keys given a value to be queried; that is, `query` is the value on the right-hand side of an indexable operator whose left-hand side is the indexed column. `n` is the strategy number of the operator within the operator class (see [Section 39.16.2](#)). Often, `extractQuery` will need to consult `n` to determine the data type of `query` and the method it should use to extract key values. The number of returned keys must be stored into `*nkeys`. If any of the keys can be null, also palloc an array of `*nkeys` `bool` fields, store its address at `*nullFlags`, and set these null flags as needed. `*nullFlags` can be left `NULL` (its initial value) if all keys are non-null. The return value can be `NULL` if the query contains no keys.

`searchMode` is an output argument that allows `extractQuery` to specify details about how the search will be done. If `*searchMode` is set to `GIN_SEARCH_MODE_DEFAULT` (which is the value it is initialized to before call), only items that match at least one of the returned keys are considered candidate matches. If `*searchMode` is set to `GIN_SEARCH_MODE_INCLUDE_EMPTY`, then in addition to items containing at least one matching key, items that contain no keys at all are considered candidate matches. (This mode is useful for implementing is-subset-of operators, for example.) If `*searchMode` is set to `GIN_SEARCH_MODE_ALL`, then all non-null items in the index are considered candidate matches, whether they match any of the returned keys or not. (This mode is much slower than the other two choices, since it requires scanning essentially the entire index, but it may be necessary to implement

corner cases correctly. An operator that needs this mode in most cases is probably not a good candidate for a GIN operator class.) The symbols to use for setting this mode are defined in `access/gin.h`.

`pmatch` is an output argument for use when partial match is supported. To use it, `extractQuery` must allocate an array of `*nkeys` `bools` and store its address at `*pmatch`. Each element of the array should be set to `true` if the corresponding key requires partial match, `false` if not. If `*pmatch` is set to `NULL` then GIN assumes partial match is not required. The variable is initialized to `NULL` before call, so this argument can simply be ignored by operator classes that do not support partial match.

`extra_data` is an output argument that allows `extractQuery` to pass additional data to the `consistent` and `comparePartial` methods. To use it, `extractQuery` must allocate an array of `*nkeys` `pointers` and store its address at `*extra_data`, then store whatever it wants to into the individual pointers. The variable is initialized to `NULL` before call, so this argument can simply be ignored by operator classes that do not require extra data. If `*extra_data` is set, the whole array is passed to the `consistent` method, and the appropriate element to the `comparePartial` method.

An operator class must also provide a function to check if an indexed item matches the query. It comes in two flavors, a Boolean `consistent` function, and a ternary `triConsistent` function. `triConsistent` covers the functionality of both, so providing `triConsistent` alone is sufficient. However, if the Boolean variant is significantly cheaper to calculate, it can be advantageous to provide both. If only the Boolean variant is provided, some optimizations that depend on refuting index items before fetching all the keys are disabled.

```
bool consistent(bool check[], StrategyNumber n, Datum query, int32 nkeys, Pointer extra_data[], bool *recheck, Datum queryKeys[], bool nullFlags[])
```

Returns `true` if an indexed item satisfies the query operator with strategy number `n` (or might satisfy it, if the `recheck` indication is returned). This function does not have direct access to the indexed item's value, since GIN does not store items explicitly. Rather, what is available is knowledge about which key values extracted from the query appear in a given indexed item. The `check` array has length `nkeys`, which is the same as the number of keys previously returned by `extractQuery` for this query datum. Each element of the `check` array is `true` if the indexed item contains the corresponding query key, i.e., if `(check[i] == true)` the `i`-th key of the `extractQuery` result array is present in the indexed item. The original query datum is passed in case the `consistent` method needs to consult it, and so are the `queryKeys[]` and `nullFlags[]` arrays previously returned by `extractQuery`. `extra_data` is the extra-data array returned by `extractQuery`, or `NULL` if none.

When `extractQuery` returns a null key in `queryKeys[]`, the corresponding `check[]` element is `true` if the indexed item contains a null key; that is, the semantics of `check[]` are like `IS NOT DISTINCT FROM`. The `consistent` function can examine the corresponding `nullFlags[]` element if it needs to tell the difference between a regular value match and a null match.

On success, `*recheck` should be set to `true` if the heap tuple needs to be rechecked against the query operator, or `false` if the index test is exact. That is, a false return value guarantees that the heap tuple does not match the query; a true return value with `*recheck` set to `false` guarantees that the heap tuple does match the query; and a true return value with `*recheck` set to `true` means that the heap tuple might match the query, so it needs to be fetched and rechecked by evaluating the query operator directly against the originally indexed item.

```
GinTernaryValue triConsistent(GinTernaryValue check[], StrategyNumber n, Datum query, int32 nkeys, Pointer extra_data[], Datum queryKeys[], bool nullFlags[])
```

`triConsistent` is similar to `consistent`, but instead of Booleans in the `check` vector, there are three possible values for each key: `GIN_TRUE`, `GIN_FALSE` and `GIN_MAYBE`. `GIN_FALSE` and `GIN_TRUE` have the same meaning as regular Boolean values, while `GIN_MAYBE` means that the presence of that key is not known. When `GIN_MAYBE` values are present, the function should only return `GIN_TRUE` if the item certainly matches whether or not the index item contains the corresponding query keys. Likewise, the function must return `GIN_FALSE` only if the item certainly does not match, whether or not it contains the `GIN_MAYBE` keys. If the result depends on the `GIN_MAYBE` entries, i.e., the match cannot be confirmed or refuted based on the known query keys, the function must return `GIN_MAYBE`.

When there are no `GIN_MAYBE` values in the `check` vector, a `GIN_MAYBE` return value is the equivalent of setting the `recheck` flag in the Boolean `consistent` function.

In addition, GIN must have a way to sort the key values stored in the index. The operator class can define the sort ordering by specifying a comparison method:

```
int compare(Datum a, Datum b)
```

Compares two keys (not indexed items!) and returns an integer less than zero, zero, or greater than zero, indicating whether the first key is less than, equal to, or greater than the second. Null keys are never passed to this function.

Alternatively, if the operator class does not provide a `compare` method, GIN will look up the default btree operator class for the index key data type, and use its comparison function. It is recommended to specify the comparison function in a GIN operator class that is meant for just one data type, as looking up the btree operator class costs a few cycles. However, polymorphic GIN operator classes (such as `array_ops`) typically cannot specify a single comparison function.

An operator class for GIN can optionally supply the following methods:

```
int comparePartial(Datum partial_key, Datum key, StrategyNumber n, Pointer extra_data)
```

Compare a partial-match query key to an index key. Returns an integer whose sign indicates the result: less than zero means the index key does not match the query, but the index scan should continue; zero means that the index key does match the query; greater than zero indicates that the index scan should stop because no more matches are possible. The strategy number `n` of the operator that generated the partial match query is provided, in case its semantics are needed to determine when to end the scan. Also, `extra_data` is the corresponding element of the extra-data array made by `extractQuery`, or `NULL` if none. Null keys are never passed to this function.

```
void options(local_relopts *relopts)
```

Defines a set of user-visible parameters that control operator class behavior.

The `options` function is passed a pointer to a `local_relopts` struct, which needs to be filled with a set of operator class specific options. The options can be accessed from other support functions using the `PG_HAS_OPCLASS_OPTIONS()` and `PG_GET_OPCLASS_OPTIONS()` macros.

Since both key extraction of indexed values and representation of the key in GIN are flexible, they may depend on user-specified parameters.

To support “partial match” queries, an operator class must provide the `comparePartial` method, and its `extractQuery` method must set the `pmatch` parameter when a partial-match query is encountered. See [Section 65.4.4.2](#) for details.

The actual data types of the various `Datum` values mentioned above vary depending on the operator class. The item values passed to `extractValue` are always of the operator class's input type, and all key values must be of the class's `STORAGE` type. The type of the `query` argument passed to `extractQuery`, `consistent` and `triConsistent` is whatever is the right-hand input type of the class member operator identified by the strategy number. This need not be the same as the indexed type, so long as key values of the correct type can be extracted from it. However, it is recommended that the SQL declarations of these three support functions use the opclass's indexed data type for the `query` argument, even though the actual type might be something else depending on the operator.

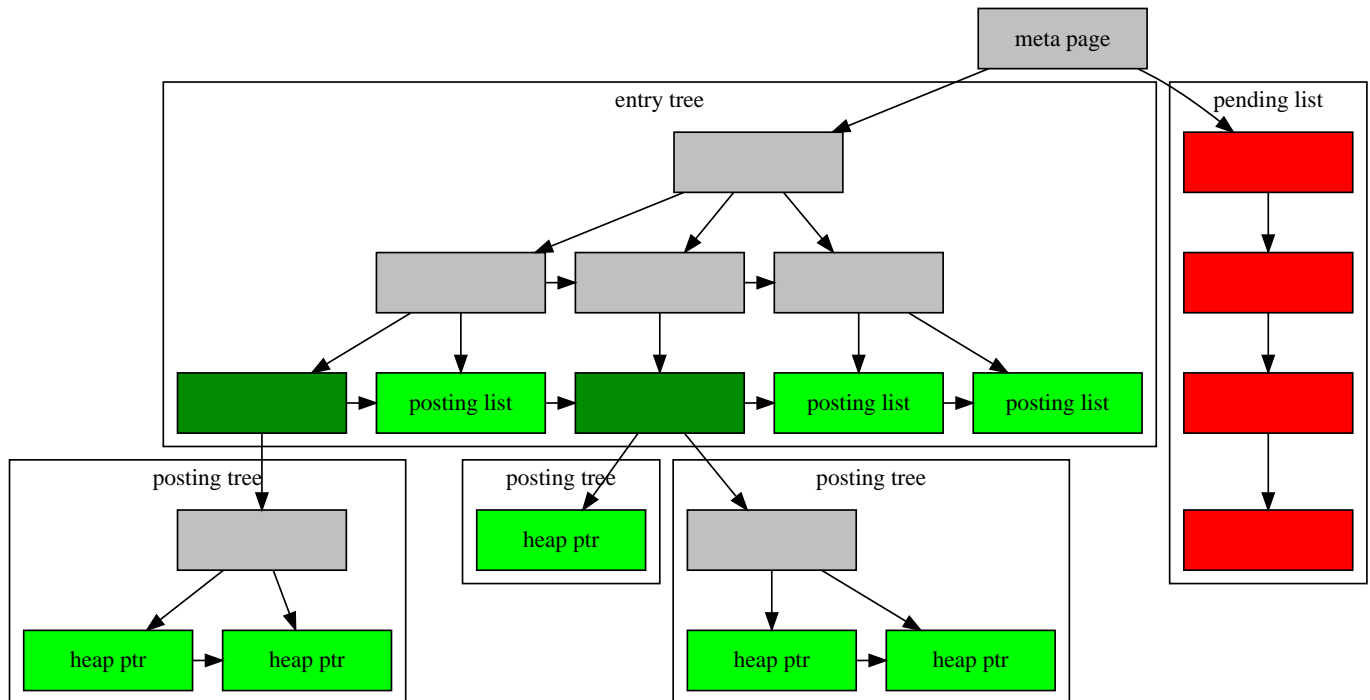
65.4.4. Implementation

Internally, a GIN index contains a B-tree index constructed over keys, where each key is an element of one or more indexed items (a member of an array, for example) and where each tuple in a leaf page contains either a pointer to a B-tree of heap pointers (a “posting tree”), or a simple list of heap pointers (a “posting list”) when the list is small enough to fit into a single index tuple along with the key value. [Figure 65.1](#) illustrates these components of a GIN index.

As of PostgreSQL 9.1, null key values can be included in the index. Also, placeholder nulls are included in the index for indexed items that are null or contain no keys according to `extractValue`. This allows searches that should find empty items to do so.

Multicolumn GIN indexes are implemented by building a single B-tree over composite values (column number, key value). The key values for different columns can be of different types.

Figure 65.1. GIN Internals



65.4.4.1. GIN Fast Update Technique

Updating a GIN index tends to be slow because of the intrinsic nature of inverted indexes: inserting or updating one heap row can cause many inserts into the index (one for each key extracted from the indexed item). GIN is capable of postponing much of this work by inserting new tuples into a temporary, unsorted list of pending entries. When the table is vacuumed or autoanalyzed, or when `gin_clean_pending_list` function is called, or if the pending list becomes larger than `gin_pending_list_limit`, the entries are moved to the main GIN data structure using the same bulk insert techniques used during initial index creation. This greatly improves GIN index update speed, even counting the additional vacuum overhead. Moreover the overhead work can be done by a background process instead of in foreground query processing.

The main disadvantage of this approach is that searches must scan the list of pending entries in addition to searching the regular index, and so a large list of pending entries will slow searches significantly. Another disadvantage is that, while most updates are fast, an update that causes the pending list to become “too large” will incur an immediate cleanup cycle and thus be much slower than other updates. Proper use of autovacuum can minimize both of these problems.

If consistent response time is more important than update speed, use of pending entries can be disabled by turning off the `fastupdate` storage parameter for a GIN index. See [CREATE INDEX](#) for details.

65.4.4.2. Partial Match Algorithm

GIN can support “partial match” queries, in which the query does not determine an exact match for one or more keys, but the possible matches fall within a reasonably narrow range of key values (within the key sorting order determined by the `compare` support method). The `extractQuery` method, instead of returning a key value to be matched exactly, returns a key value that is the lower bound of the range to be searched, and sets the `pmatch` flag true. The key range is then scanned using the `comparePartial`

method. `comparePartial` must return zero for a matching index key, less than zero for a non-match that is still within the range to be searched, or greater than zero if the index key is past the range that could match.

65.4.5. GIN Tips and Tricks

Create vs. insert

Insertion into a GIN index can be slow due to the likelihood of many keys being inserted for each item. So, for bulk insertions into a table it is advisable to drop the GIN index and recreate it after finishing bulk insertion.

When `fastupdate` is enabled for GIN (see [Section 65.4.4.1](#) for details), the penalty is less than when it is not. But for very large updates it may still be best to drop and recreate the index.

`maintenance_work_mem`

Build time for a GIN index is very sensitive to the `maintenance_work_mem` setting; it doesn't pay to skimp on work memory during index creation.

`gin_pending_list_limit`

During a series of insertions into an existing GIN index that has `fastupdate` enabled, the system will clean up the pending-entry list whenever the list grows larger than `gin_pending_list_limit`. To avoid fluctuations in observed response time, it's desirable to have pending-list cleanup occur in the background (i.e., via autovacuum). Foreground cleanup operations can be avoided by increasing `gin_pending_list_limit` or making autovacuum more aggressive. However, enlarging the threshold of the cleanup operation means that if a foreground cleanup does occur, it will take even longer.

`gin_pending_list_limit` can be overridden for individual GIN indexes by changing storage parameters, which allows each GIN index to have its own cleanup threshold. For example, it's possible to increase the threshold only for the GIN index which can be updated heavily, and decrease it otherwise.

`gin_fuzzy_search_limit`

The primary goal of developing GIN indexes was to create support for highly scalable full-text search in Postgres Pro, and there are often situations when a full-text search returns a very large set of results. Moreover, this often happens when the query contains very frequent words, so that the large result set is not even useful. Since reading many tuples from the disk and sorting them could take a lot of time, this is unacceptable for production. (Note that the index search itself is very fast.)

To facilitate controlled execution of such queries, GIN has a configurable soft upper limit on the number of rows returned: the `gin_fuzzy_search_limit` configuration parameter. It is set to 0 (meaning no limit) by default. If a non-zero limit is set, then the returned set is a subset of the whole result set, chosen at random.

“Soft” means that the actual number of returned results could differ somewhat from the specified limit, depending on the query and the quality of the system's random number generator.

From experience, values in the thousands (e.g., 5000 — 20000) work well.

65.4.6. Limitations

GIN assumes that indexable operators are strict. This means that `extractValue` will not be called at all on a null item value (instead, a placeholder index entry is created automatically), and `extractQuery` will not be called on a null query value either (instead, the query is presumed to be unsatisfiable). Note however that null key values contained within a non-null composite item or query value are supported.

65.4.7. Examples

The core Postgres Pro distribution includes the GIN operator classes previously shown in [Table 65.3](#). The following `contrib` modules also contain GIN operator classes:

`btree_gin`

B-tree equivalent functionality for several data types

`hstore`

Module for storing (key, value) pairs

`intarray`

Enhanced support for `int[]`

`pg_trgm`

Text similarity using trigram matching

65.5. BRIN Indexes

65.5.1. Introduction

BRIN stands for Block Range Index. BRIN is designed for handling very large tables in which certain columns have some natural correlation with their physical location within the table.

BRIN works in terms of *block ranges* (or “page ranges”). A block range is a group of pages that are physically adjacent in the table; for each block range, some summary info is stored by the index. For example, a table storing a store's sale orders might have a date column on which each order was placed, and most of the time the entries for earlier orders will appear earlier in the table as well; a table storing a ZIP code column might have all codes for a city grouped together naturally.

BRIN indexes can satisfy queries via regular bitmap index scans, and will return all tuples in all pages within each range if the summary info stored by the index is *consistent* with the query conditions. The query executor is in charge of rechecking these tuples and discarding those that do not match the query conditions — in other words, these indexes are lossy. Because a BRIN index is very small, scanning the index adds little overhead compared to a sequential scan, but may avoid scanning large parts of the table that are known not to contain matching tuples.

The specific data that a BRIN index will store, as well as the specific queries that the index will be able to satisfy, depend on the operator class selected for each column of the index. Data types having a linear sort order can have operator classes that store the minimum and maximum value within each block range, for instance; geometrical types might store the bounding box for all the objects in the block range.

The size of the block range is determined at index creation time by the `pages_per_range` storage parameter. The number of index entries will be equal to the size of the relation in pages divided by the selected value for `pages_per_range`. Therefore, the smaller the number, the larger the index becomes (because of the need to store more index entries), but at the same time the summary data stored can be more precise and more data blocks can be skipped during an index scan.

65.5.1.1. Index Maintenance

At the time of creation, all existing heap pages are scanned and a summary index tuple is created for each range, including the possibly-incomplete range at the end. As new pages are filled with data, page ranges that are already summarized will cause the summary information to be updated with data from the new tuples. When a new page is created that does not fall within the last summarized range, the range that the new page belongs to does not automatically acquire a summary tuple; those tuples remain unsummarized until a summarization run is invoked later, creating the initial summary for that range.

There are several ways to trigger the initial summarization of a page range. If the table is vacuumed, either manually or by [autovacuum](#), all existing unsummarized page ranges are summarized. Also, if the index's [autosummarize](#) parameter is enabled, which it isn't by default, whenever [autovacuum](#) runs in that database, summarization will occur for all unsummarized page ranges that have been filled, regardless of whether the table itself is processed by [autovacuum](#); see below.

Lastly, the following functions can be used (while these functions run, [search_path](#) is temporarily changed to `pg_catalog`, `pg_temp`):

`brin_summarize_new_values(regclass)` which summarizes all unsummarized ranges;
`brin_summarize_range(regclass, bigint)` which summarizes only the range containing the given page, if it is unsummarized.

When autosummarization is enabled, a request is sent to `autovacuum` to execute a targeted summarization for a block range when an insertion is detected for the first item of the first page of the next block range, to be fulfilled the next time an `autovacuum` worker finishes running in the same database. If the request queue is full, the request is not recorded and a message is sent to the server log:

```
LOG:  request for BRIN range summarization for index "brin_wi_idx" page 128 was not
recorded
```

When this happens, the range will remain unsummarized until the next regular vacuum run on the table, or one of the functions mentioned above are invoked.

Conversely, a range can be de-summarized using the `brin_desummarize_range(regclass, bigint)` function, which is useful when the index tuple is no longer a very good representation because the existing values have changed. See [Section 9.28.8](#) for details.

65.5.2. Built-in Operator Classes

The core Postgres Pro distribution includes the BRIN operator classes shown in [Table 65.4](#).

The *minmax* operator classes store the minimum and the maximum values appearing in the indexed column within the range. The *inclusion* operator classes store a value which includes the values in the indexed column within the range. The *bloom* operator classes build a Bloom filter for all values in the range. The *minmax-multi* operator classes store multiple minimum and maximum values, representing values appearing in the indexed column within the range.

Table 65.4. Built-in BRIN Operator Classes

Name	Indexable Operators
bit_minmax_ops	= (bit,bit)
	< (bit,bit)
	> (bit,bit)
	<= (bit,bit)
	>= (bit,bit)
box_inclusion_ops	@> (box,point)
	<< (box,box)
	&< (box,box)
	&> (box,box)
	>> (box,box)
	<@ (box,box)
	@> (box,box)
	~= (box,box)
	&& (box,box)
	<< (box,box)
	&< (box,box)
	&> (box,box)
	>> (box,box)
bpchar_bloom_ops	= (character,character)
bpchar_minmax_ops	= (character,character)
	< (character,character)

Name	Indexable Operators
	<= (character, character)
	> (character, character)
	>= (character, character)
bytea_bloom_ops	= (bytea, bytea)
bytea_minmax_ops	= (bytea, bytea)
	< (bytea, bytea)
	<= (bytea, bytea)
	> (bytea, bytea)
	>= (bytea, bytea)
char_bloom_ops	= ("char", "char")
char_minmax_ops	= ("char", "char")
	< ("char", "char")
	<= ("char", "char")
	> ("char", "char")
	>= ("char", "char")
date_bloom_ops	= (date, date)
date_minmax_ops	= (date, date)
	< (date, date)
	<= (date, date)
	> (date, date)
	>= (date, date)
date_minmax_multi_ops	= (date, date)
	< (date, date)
	<= (date, date)
	> (date, date)
	>= (date, date)
float4_bloom_ops	= (float4, float4)
float4_minmax_ops	= (float4, float4)
	< (float4, float4)
	> (float4, float4)
	<= (float4, float4)
	>= (float4, float4)
float4_minmax_multi_ops	= (float4, float4)
	< (float4, float4)
	> (float4, float4)
	<= (float4, float4)
	>= (float4, float4)
float8_bloom_ops	= (float8, float8)
float8_minmax_ops	= (float8, float8)
	< (float8, float8)

Name	Indexable Operators
	<= (float8,float8)
	> (float8,float8)
	>= (float8,float8)
float8_minmax_multi_ops	= (float8,float8)
	< (float8,float8)
	<= (float8,float8)
	> (float8,float8)
	>= (float8,float8)
inet_inclusion_ops	<< (inet,inet)
	<<= (inet,inet)
	>> (inet,inet)
	>>= (inet,inet)
	= (inet,inet)
	&& (inet,inet)
inet_bloom_ops	= (inet,inet)
inet_minmax_ops	= (inet,inet)
	< (inet,inet)
	<= (inet,inet)
	> (inet,inet)
	>= (inet,inet)
inet_minmax_multi_ops	= (inet,inet)
	< (inet,inet)
	<= (inet,inet)
	> (inet,inet)
	>= (inet,inet)
int2_bloom_ops	= (int2,int2)
int2_minmax_ops	= (int2,int2)
	< (int2,int2)
	> (int2,int2)
	<= (int2,int2)
	>= (int2,int2)
int2_minmax_multi_ops	= (int2,int2)
	< (int2,int2)
	> (int2,int2)
	<= (int2,int2)
	>= (int2,int2)
int4_bloom_ops	= (int4,int4)
int4_minmax_ops	= (int4,int4)
	< (int4,int4)
	> (int4,int4)

Name	Indexable Operators
	<= (int4,int4)
	>= (int4,int4)
int4_minmax_multi_ops	= (int4,int4)
	< (int4,int4)
	> (int4,int4)
	<= (int4,int4)
	>= (int4,int4)
int8_bloom_ops	= (bigint,bigint)
int8_minmax_ops	= (bigint,bigint)
	< (bigint,bigint)
	> (bigint,bigint)
	<= (bigint,bigint)
	>= (bigint,bigint)
int8_minmax_multi_ops	= (bigint,bigint)
	< (bigint,bigint)
	> (bigint,bigint)
	<= (bigint,bigint)
	>= (bigint,bigint)
interval_bloom_ops	= (interval,interval)
interval_minmax_ops	= (interval,interval)
	< (interval,interval)
	<= (interval,interval)
	> (interval,interval)
	>= (interval,interval)
interval_minmax_multi_ops	= (interval,interval)
	< (interval,interval)
	<= (interval,interval)
	> (interval,interval)
	>= (interval,interval)
macaddr_bloom_ops	= (macaddr,macaddr)
macaddr_minmax_ops	= (macaddr,macaddr)
	< (macaddr,macaddr)
	<= (macaddr,macaddr)
	> (macaddr,macaddr)
	>= (macaddr,macaddr)
macaddr_minmax_multi_ops	= (macaddr,macaddr)
	< (macaddr,macaddr)
	<= (macaddr,macaddr)
	> (macaddr,macaddr)
	>= (macaddr,macaddr)

Name	Indexable Operators
macaddr8_bloom_ops	= (macaddr8,macaddr8)
macaddr8_minmax_ops	= (macaddr8,macaddr8)
	< (macaddr8,macaddr8)
	<= (macaddr8,macaddr8)
	> (macaddr8,macaddr8)
	>= (macaddr8,macaddr8)
macaddr8_minmax_multi_ops	= (macaddr8,macaddr8)
	< (macaddr8,macaddr8)
	<= (macaddr8,macaddr8)
	> (macaddr8,macaddr8)
	>= (macaddr8,macaddr8)
name_bloom_ops	= (name,name)
name_minmax_ops	= (name,name)
	< (name,name)
	<= (name,name)
	> (name,name)
	>= (name,name)
numeric_bloom_ops	= (numeric,numeric)
numeric_minmax_ops	= (numeric,numeric)
	< (numeric,numeric)
	<= (numeric,numeric)
	> (numeric,numeric)
	>= (numeric,numeric)
numeric_minmax_multi_ops	= (numeric,numeric)
	< (numeric,numeric)
	<= (numeric,numeric)
	> (numeric,numeric)
	>= (numeric,numeric)
oid_bloom_ops	= (oid,oid)
oid_minmax_ops	= (oid,oid)
	< (oid,oid)
	> (oid,oid)
	<= (oid,oid)
	>= (oid,oid)
oid_minmax_multi_ops	= (oid,oid)
	< (oid,oid)
	> (oid,oid)
	<= (oid,oid)
	>= (oid,oid)
pg_lsn_bloom_ops	= (pg_lsn,pg_lsn)

Name	Indexable Operators
pg_lsn_minmax_ops	= (pg_lsn,pg_lsn)
	< (pg_lsn,pg_lsn)
	> (pg_lsn,pg_lsn)
	<= (pg_lsn,pg_lsn)
	>= (pg_lsn,pg_lsn)
pg_lsn_minmax_multi_ops	= (pg_lsn,pg_lsn)
	< (pg_lsn,pg_lsn)
	> (pg_lsn,pg_lsn)
	<= (pg_lsn,pg_lsn)
	>= (pg_lsn,pg_lsn)
range_inclusion_ops	= (anyrange,anyrange)
	< (anyrange,anyrange)
	<= (anyrange,anyrange)
	>= (anyrange,anyrange)
	> (anyrange,anyrange)
	&& (anyrange,anyrange)
	@> (anyrange,anyelement)
	@> (anyrange,anyrange)
	<@ (anyrange,anyrange)
	<< (anyrange,anyrange)
	>> (anyrange,anyrange)
	&< (anyrange,anyrange)
	&> (anyrange,anyrange)
	- - (anyrange,anyrange)
text_bloom_ops	= (text,text)
text_minmax_ops	= (text,text)
	< (text,text)
	<= (text,text)
	> (text,text)
	>= (text,text)
tid_bloom_ops	= (tid,tid)
tid_minmax_ops	= (tid,tid)
	< (tid,tid)
	> (tid,tid)
	<= (tid,tid)
	>= (tid,tid)
tid_minmax_multi_ops	= (tid,tid)
	< (tid,tid)
	> (tid,tid)
	<= (tid,tid)

Name	Indexable Operators
	>= (tid,tid)
timestamp_bloom_ops	= (timestamp,timestamp)
timestamp_minmax_ops	= (timestamp,timestamp)
	< (timestamp,timestamp)
	<= (timestamp,timestamp)
	> (timestamp,timestamp)
	>= (timestamp,timestamp)
timestamp_minmax_multi_ops	= (timestamp,timestamp)
	< (timestamp,timestamp)
	<= (timestamp,timestamp)
	> (timestamp,timestamp)
	>= (timestamp,timestamp)
timestampptz_bloom_ops	= (timestampptz,timestampptz)
timestampptz_minmax_ops	= (timestampptz,timestampptz)
	< (timestampptz,timestampptz)
	<= (timestampptz,timestampptz)
	> (timestampptz,timestampptz)
	>= (timestampptz,timestampptz)
timestampptz_minmax_multi_ops	= (timestampptz,timestampptz)
	< (timestampptz,timestampptz)
	<= (timestampptz,timestampptz)
	> (timestampptz,timestampptz)
	>= (timestampptz,timestampptz)
time_bloom_ops	= (time,time)
time_minmax_ops	= (time,time)
	< (time,time)
	<= (time,time)
	> (time,time)
	>= (time,time)
time_minmax_multi_ops	= (time,time)
	< (time,time)
	<= (time,time)
	> (time,time)
	>= (time,time)
timetz_bloom_ops	= (timetz,timetz)
timetz_minmax_ops	= (timetz,timetz)
	< (timetz,timetz)
	<= (timetz,timetz)
	> (timetz,timetz)
	>= (timetz,timetz)

Name	Indexable Operators
timetz_minmax_multi_ops	= (timetz,timetz)
	< (timetz,timetz)
	<= (timetz,timetz)
	> (timetz,timetz)
	>= (timetz,timetz)
uuid_bloom_ops	= (uuid,uuid)
uuid_minmax_ops	= (uuid,uuid)
	< (uuid,uuid)
	> (uuid,uuid)
	<= (uuid,uuid)
	>= (uuid,uuid)
uuid_minmax_multi_ops	= (uuid,uuid)
	< (uuid,uuid)
	> (uuid,uuid)
	<= (uuid,uuid)
	>= (uuid,uuid)
varbit_minmax_ops	= (varbit,varbit)
	< (varbit,varbit)
	> (varbit,varbit)
	<= (varbit,varbit)
	>= (varbit,varbit)

65.5.2.1. Operator Class Parameters

Some of the built-in operator classes allow specifying parameters affecting behavior of the operator class. Each operator class has its own set of allowed parameters. Only the `bloom` and `minmax-multi` operator classes allow specifying parameters:

`bloom` operator classes accept these parameters:

`n_distinct_per_range`

Defines the estimated number of distinct non-null values in the block range, used by BRIN bloom indexes for sizing of the Bloom filter. It behaves similarly to `n_distinct` option for [ALTER TABLE](#). When set to a positive value, each block range is assumed to contain this number of distinct non-null values. When set to a negative value, which must be greater than or equal to -1, the number of distinct non-null values is assumed to grow linearly with the maximum possible number of tuples in the block range (about 290 rows per block). The default value is `-0.1`, and the minimum number of distinct non-null values is 16.

`false_positive_rate`

Defines the desired false positive rate used by BRIN bloom indexes for sizing of the Bloom filter. The values must be between 0.0001 and 0.25. The default value is 0.01, which is 1% false positive rate.

`minmax-multi` operator classes accept these parameters:

`values_per_range`

Defines the maximum number of values stored by BRIN minmax indexes to summarize a block range. Each value may represent either a point, or a boundary of an interval. Values must be between 8 and 256, and the default value is 32.

65.5.3. Extensibility

The BRIN interface has a high level of abstraction, requiring the access method implementer only to implement the semantics of the data type being accessed. The BRIN layer itself takes care of concurrency, logging and searching the index structure.

All it takes to get a BRIN access method working is to implement a few user-defined methods, which define the behavior of summary values stored in the index and the way they interact with scan keys. In short, BRIN combines extensibility with generality, code reuse, and a clean interface.

There are four methods that an operator class for BRIN must provide:

```
BrinOpcInfo *opcInfo(Oid type_oid)
```

Returns internal information about the indexed columns' summary data. The return value must point to a `palloc'd` `BrinOpcInfo`, which has this definition:

```
typedef struct BrinOpcInfo
{
    /* Number of columns stored in an index column of this opclass */
    uint16      oi_nstored;

    /* Opaque pointer for the opclass' private use */
    void        *oi_opaque;

    /* Type cache entries of the stored columns */
    TypeCacheEntry *oi_tpcache[FLEXIBLE_ARRAY_MEMBER];
} BrinOpcInfo;
```

`BrinOpcInfo.oi_opaque` can be used by the operator class routines to pass information between support functions during an index scan.

```
bool consistent(BrinDesc *bdesc, BrinValues *column, ScanKey *keys, int nkeys)
```

Returns whether all the `ScanKey` entries are consistent with the given indexed values for a range. The attribute number to use is passed as part of the scan key. Multiple scan keys for the same attribute may be passed at once; the number of entries is determined by the `nkeys` parameter.

```
bool consistent(BrinDesc *bdesc, BrinValues *column, ScanKey key)
```

Returns whether the `ScanKey` is consistent with the given indexed values for a range. The attribute number to use is passed as part of the scan key. This is an older backward-compatible variant of the `consistent` function.

```
bool addValue(BrinDesc *bdesc, BrinValues *column, Datum newval, bool isnull)
```

Given an index tuple and an indexed value, modifies the indicated attribute of the tuple so that it additionally represents the new value. If any modification was done to the tuple, `true` is returned.

```
bool unionTuples(BrinDesc *bdesc, BrinValues *a, BrinValues *b)
```

Consolidates two index tuples. Given two index tuples, modifies the indicated attribute of the first of them so that it represents both tuples. The second tuple is not modified.

An operator class for BRIN can optionally specify the following method:

```
void options(local_relopts *relopts)
```

Defines a set of user-visible parameters that control operator class behavior.

The `options` function is passed a pointer to a `local_relopts` struct, which needs to be filled with a set of operator class specific options. The options can be accessed from other support functions using the `PG_HAS_OPCLASS_OPTIONS()` and `PG_GET_OPCLASS_OPTIONS()` macros.

Since both key extraction of indexed values and representation of the key in BRIN are flexible, they may depend on user-specified parameters.

The core distribution includes support for four types of operator classes: minmax, minmax-multi, inclusion and bloom. Operator class definitions using them are shipped for in-core data types as appropriate. Additional operator classes can be defined by the user for other data types using equivalent definitions, without having to write any source code; appropriate catalog entries being declared is enough. Note that assumptions about the semantics of operator strategies are embedded in the support functions' source code.

Operator classes that implement completely different semantics are also possible, provided implementations of the four main support functions described above are written. Note that backwards compatibility across major releases is not guaranteed: for example, additional support functions might be required in later releases.

To write an operator class for a data type that implements a totally ordered set, it is possible to use the minmax support functions alongside the corresponding operators, as shown in [Table 65.5](#). All operator class members (functions and operators) are mandatory.

Table 65.5. Function and Support Numbers for Minmax Operator Classes

Operator class member	Object
Support Function 1	internal function <code>brin_minmax_opcinfo()</code>
Support Function 2	internal function <code>brin_minmax_add_value()</code>
Support Function 3	internal function <code>brin_minmax_consistent()</code>
Support Function 4	internal function <code>brin_minmax_union()</code>
Operator Strategy 1	operator less-than
Operator Strategy 2	operator less-than-or-equal-to
Operator Strategy 3	operator equal-to
Operator Strategy 4	operator greater-than-or-equal-to
Operator Strategy 5	operator greater-than

To write an operator class for a complex data type which has values included within another type, it's possible to use the inclusion support functions alongside the corresponding operators, as shown in [Table 65.6](#). It requires only a single additional function, which can be written in any language. More functions can be defined for additional functionality. All operators are optional. Some operators require other operators, as shown as dependencies on the table.

Table 65.6. Function and Support Numbers for Inclusion Operator Classes

Operator class member	Object	Dependency
Support Function 1	internal function <code>brin_inclusion_opcinfo()</code>	
Support Function 2	internal function <code>brin_inclusion_add_value()</code>	
Support Function 3	internal function <code>brin_inclusion_consistent()</code>	
Support Function 4	internal function <code>brin_inclusion_union()</code>	
Support Function 11	function to merge two elements	
Support Function 12	optional function to check whether two elements are mergeable	
Support Function 13	optional function to check if an element is contained within another	

Operator class member	Object	Dependency
Support Function 14	optional function to check whether an element is empty	
Operator Strategy 1	operator left-of	Operator Strategy 4
Operator Strategy 2	operator does-not-extend-to-the-right-of	Operator Strategy 5
Operator Strategy 3	operator overlaps	
Operator Strategy 4	operator does-not-extend-to-the-left-of	Operator Strategy 1
Operator Strategy 5	operator right-of	Operator Strategy 2
Operator Strategy 6, 18	operator same-as-or-equal-to	Operator Strategy 7
Operator Strategy 7, 16, 24, 25	operator contains-or-equal-to	
Operator Strategy 8, 26, 27	operator is-contained-by-or-equal-to	Operator Strategy 3
Operator Strategy 9	operator does-not-extend-above	Operator Strategy 11
Operator Strategy 10	operator is-below	Operator Strategy 12
Operator Strategy 11	operator is-above	Operator Strategy 9
Operator Strategy 12	operator does-not-extend-below	Operator Strategy 10
Operator Strategy 20	operator less-than	Operator Strategy 5
Operator Strategy 21	operator less-than-or-equal-to	Operator Strategy 5
Operator Strategy 22	operator greater-than	Operator Strategy 1
Operator Strategy 23	operator greater-than-or-equal-to	Operator Strategy 1

Support function numbers 1 through 10 are reserved for the BRIN internal functions, so the SQL level functions start with number 11. Support function number 11 is the main function required to build the index. It should accept two arguments with the same data type as the operator class, and return the union of them. The inclusion operator class can store union values with different data types if it is defined with the `STORAGE` parameter. The return value of the union function should match the `STORAGE` data type.

Support function numbers 12 and 14 are provided to support irregularities of built-in data types. Function number 12 is used to support network addresses from different families which are not mergeable. Function number 14 is used to support empty ranges. Function number 13 is an optional but recommended one, which allows the new value to be checked before it is passed to the union function. As the BRIN framework can shortcut some operations when the union is not changed, using this function can improve index performance.

To write an operator class for a data type that implements only an equality operator and supports hashing, it is possible to use the bloom support procedures alongside the corresponding operators, as shown in [Table 65.7](#). All operator class members (procedures and operators) are mandatory.

Table 65.7. Procedure and Support Numbers for Bloom Operator Classes

Operator class member	Object
Support Procedure 1	internal function <code>brin_bloom_opcinfo()</code>
Support Procedure 2	internal function <code>brin_bloom_add_value()</code>
Support Procedure 3	internal function <code>brin_bloom_consistent()</code>
Support Procedure 4	internal function <code>brin_bloom_union()</code>
Support Procedure 5	internal function <code>brin_bloom_options()</code>
Support Procedure 11	function to compute hash of an element
Operator Strategy 1	operator equal-to

Support procedure numbers 1-10 are reserved for the BRIN internal functions, so the SQL level functions start with number 11. Support function number 11 is the main function required to build the index. It should accept one argument with the same data type as the operator class, and return a hash of the value.

The minmax-multi operator class is also intended for data types implementing a totally ordered set, and may be seen as a simple extension of the minmax operator class. While minmax operator class summarizes values from each block range into a single contiguous interval, minmax-multi allows summarization into multiple smaller intervals to improve handling of outlier values. It is possible to use the minmax-multi support procedures alongside the corresponding operators, as shown in [Table 65.8](#). All operator class members (procedures and operators) are mandatory.

Table 65.8. Procedure and Support Numbers for minmax-multi Operator Classes

Operator class member	Object
Support Procedure 1	internal function <code>brin_minmax_multi_opcinfo()</code>
Support Procedure 2	internal function <code>brin_minmax_multi_add_value()</code>
Support Procedure 3	internal function <code>brin_minmax_multi_consistent()</code>
Support Procedure 4	internal function <code>brin_minmax_multi_union()</code>
Support Procedure 5	internal function <code>brin_minmax_multi_options()</code>
Support Procedure 11	function to compute distance between two values (length of a range)
Operator Strategy 1	operator less-than
Operator Strategy 2	operator less-than-or-equal-to
Operator Strategy 3	operator equal-to
Operator Strategy 4	operator greater-than-or-equal-to
Operator Strategy 5	operator greater-than

Both minmax and inclusion operator classes support cross-data-type operators, though with these the dependencies become more complicated. The minmax operator class requires a full set of operators to be defined with both arguments having the same data type. It allows additional data types to be supported by defining extra sets of operators. Inclusion operator class operator strategies are dependent on another operator strategy as shown in [Table 65.6](#), or the same operator strategy as themselves. They require the dependency operator to be defined with the `STORAGE` data type as the left-hand-side argument and the other supported data type to be the right-hand-side argument of the supported operator. See `float4_minmax_ops` as an example of minmax, and `box_inclusion_ops` as an example of inclusion.

65.6. Hash Indexes

65.6.1. Overview

Postgres Pro includes an implementation of persistent on-disk hash indexes, which are fully crash recoverable. Any data type can be indexed by a hash index, including data types that do not have a well-defined linear ordering. Hash indexes store only the hash value of the data being indexed, thus there are no restrictions on the size of the data column being indexed.

Hash indexes support only single-column indexes and do not allow uniqueness checking.

Hash indexes support only the `=` operator, so `WHERE` clauses that specify range operations will not be able to take advantage of hash indexes.

Each hash index tuple stores just the 4-byte hash value, not the actual column value. As a result, hash indexes may be much smaller than B-trees when indexing longer data items such as UUIDs, URLs, etc. The absence of the column value also makes all hash index scans lossy. Hash indexes may take part in bitmap index scans and backward scans.

Hash indexes are best optimized for SELECT and UPDATE-heavy workloads that use equality scans on larger tables. In a B-tree index, searches must descend through the tree until the leaf page is found. In tables with millions of rows, this descent can increase access time to data. The equivalent of a leaf page in a hash index is referred to as a bucket page. In contrast, a hash index allows accessing the bucket pages directly, thereby potentially reducing index access time in larger tables. This reduction in "logical I/O" becomes even more pronounced on indexes/data larger than shared_buffers/RAM.

Hash indexes have been designed to cope with uneven distributions of hash values. Direct access to the bucket pages works well if the hash values are evenly distributed. When inserts mean that the bucket page becomes full, additional overflow pages are chained to that specific bucket page, locally expanding the storage for index tuples that match that hash value. When scanning a hash bucket during queries, we need to scan through all of the overflow pages. Thus an unbalanced hash index might actually be worse than a B-tree in terms of number of block accesses required, for some data.

As a result of the overflow cases, we can say that hash indexes are most suitable for unique, nearly unique data or data with a low number of rows per hash bucket. One possible way to avoid problems is to exclude highly non-unique values from the index using a partial index condition, but this may not be suitable in many cases.

Like B-Trees, hash indexes perform simple index tuple deletion. This is a deferred maintenance operation that deletes index tuples that are known to be safe to delete (those whose item identifier's LP_DEAD bit is already set). If an insert finds no space is available on a page we try to avoid creating a new overflow page by attempting to remove dead index tuples. Removal cannot occur if the page is pinned at that time. Deletion of dead index pointers also occurs during VACUUM.

If it can, VACUUM will also try to squeeze the index tuples onto as few overflow pages as possible, minimizing the overflow chain. If an overflow page becomes empty, overflow pages can be recycled for reuse in other buckets, though we never return them to the operating system. There is currently no provision to shrink a hash index, other than by rebuilding it with REINDEX. There is no provision for reducing the number of buckets, either.

Hash indexes may expand the number of bucket pages as the number of rows indexed grows. The hash key-to-bucket-number mapping is chosen so that the index can be incrementally expanded. When a new bucket is to be added to the index, exactly one existing bucket will need to be "split", with some of its tuples being transferred to the new bucket according to the updated key-to-bucket-number mapping.

The expansion occurs in the foreground, which could increase execution time for user inserts. Thus, hash indexes may not be suitable for tables with rapidly increasing number of rows.

65.6.2. Implementation

There are four kinds of pages in a hash index: the meta page (page zero), which contains statically allocated control information; primary bucket pages; overflow pages; and bitmap pages, which keep track of overflow pages that have been freed and are available for re-use. For addressing purposes, bitmap pages are regarded as a subset of the overflow pages.

Both scanning the index and inserting tuples require locating the bucket where a given tuple ought to be located. To do this, we need the bucket count, highmask, and lowmask from the metapage; however, it's undesirable for performance reasons to have to lock and pin the metapage for every such operation. Instead, we retain a cached copy of the metapage in each backend's relcache entry. This will produce the correct bucket mapping as long as the target bucket hasn't been split since the last cache refresh.

Primary bucket pages and overflow pages are allocated independently since any given index might need more or fewer overflow pages relative to its number of buckets. The hash code uses an interesting set

of addressing rules to support a variable number of overflow pages while not having to move primary bucket pages around after they are created.

Each row in the table indexed is represented by a single index tuple in the hash index. Hash index tuples are stored in bucket pages, and if they exist, overflow pages. We speed up searches by keeping the index entries in any one index page sorted by hash code, thus allowing binary search to be used within an index page. Note however that there is **no** assumption about the relative ordering of hash codes across different index pages of a bucket.

The bucket splitting algorithms to expand the hash index are too complex to be worthy of mention here. The split algorithm is crash safe and can be restarted if not completed successfully.

Chapter 66. Database Physical Storage

This chapter provides an overview of the physical storage format used by Postgres Pro databases.

66.1. Database File Layout

This section describes the storage format at the level of files and directories.

Traditionally, the configuration and data files used by a database cluster are stored together within the cluster's data directory, commonly referred to as `PGDATA` (after the name of the environment variable that can be used to define it). A common location for `PGDATA` is `/var/lib/pgpro/sdm-17/data`. Multiple clusters, managed by different server instances, can exist on the same machine.

The `PGDATA` directory contains several subdirectories and control files, as shown in [Table 66.1](#). In addition to these required items, the cluster configuration files `postgresql.conf`, `pg_hba.conf`, and `pg_ident.conf` are traditionally stored in `PGDATA`, although it is possible to place them elsewhere.

Table 66.1. Contents of `PGDATA`

Item	Description
<code>PG_VERSION</code>	A file containing the major version number of Postgres Pro
<code>base</code>	Subdirectory containing per-database subdirectories
<code>current_logfiles</code>	File recording the log file(s) currently written to by the logging collector
<code>global</code>	Subdirectory containing cluster-wide tables, such as <code>pg_database</code>
<code>pg_commit_ts</code>	Subdirectory containing transaction commit timestamp data
<code>pg_dynshmem</code>	Subdirectory containing files used by the dynamic shared memory subsystem
<code>pg_logical</code>	Subdirectory containing status data for logical decoding
<code>pg_multixact</code>	Subdirectory containing multitransaction status data (used for shared row locks)
<code>pg_notify</code>	Subdirectory containing LISTEN/NOTIFY status data
<code>pg_replslot</code>	Subdirectory containing replication slot data
<code>pg_serial</code>	Subdirectory containing information about committed serializable transactions
<code>pg_snapshots</code>	Subdirectory containing exported snapshots
<code>pg_stat</code>	Subdirectory containing permanent files for the statistics subsystem
<code>pg_stat_tmp</code>	Subdirectory containing temporary files for the statistics subsystem
<code>pg_subtrans</code>	Subdirectory containing subtransaction status data
<code>pg_tblspc</code>	Subdirectory containing symbolic links to tablespaces
<code>pg_twophase</code>	Subdirectory containing state files for prepared transactions

Item	Description
<code>pg_wal</code>	Subdirectory containing WAL (Write Ahead Log) files
<code>pg_xact</code>	Subdirectory containing transaction commit status data
<code>postgresql.auto.conf</code>	A file used for storing configuration parameters that are set by <code>ALTER SYSTEM</code>
<code>postmaster.opts</code>	A file recording the command-line options the server was last started with
<code>postmaster.pid</code>	A lock file recording the current postmaster process ID (PID), cluster data directory path, postmaster start timestamp, port number, Unix-domain socket directory path (could be empty), first valid listen_address (IP address or *, or empty if not listening on TCP), and shared memory segment ID (this file is not present after server shutdown)

For each database in the cluster there is a subdirectory within `PGDATA/base`, named after the database's OID in `pg_database`. This subdirectory is the default location for the database's files; in particular, its system catalogs are stored there.

Note that the following sections describe the behavior of the builtin heap [table access method](#), and the builtin [index access methods](#). Due to the extensible nature of Postgres Pro, other access methods might work differently.

Each table and index is stored in a separate file. For ordinary relations, these files are named after the table or index's *filenode* number, which can be found in `pg_class.relfilenode`. But for temporary relations, the file name is of the form `tBBB_FFF`, where *BBB* is the process number of the backend which created the file, and *FFF* is the filenode number. In either case, in addition to the main file (a/k/a main fork), each table and index has a *free space map* (see [Section 66.3](#)), which stores information about free space available in the relation. The free space map is stored in a file named with the filenode number plus the suffix `_fsm`. Tables also have a *visibility map*, stored in a fork with the suffix `_vm`, to track which pages are known to have no dead tuples. The visibility map is described further in [Section 66.4](#). Unlogged tables and indexes have a third fork, known as the initialization fork, which is stored in a fork with the suffix `_init` (see [Section 66.5](#)).

Caution

Note that while a table's filenode often matches its OID, this is *not* necessarily the case; some operations, like `TRUNCATE`, `REINDEX`, `CLUSTER` and some forms of `ALTER TABLE`, can change the filenode while preserving the OID. Avoid assuming that filenode and table OID are the same. Also, for certain system catalogs including `pg_class` itself, `pg_class.relfilenode` contains zero. The actual filenode number of these catalogs is stored in a lower-level data structure, and can be obtained using the `pg_relation_filenode()` function.

When a table or index exceeds 1 GB, it is divided into gigabyte-sized *segments*. The first segment's file name is the same as the filenode; subsequent segments are named `filenode.1`, `filenode.2`, etc. This arrangement avoids problems on platforms that have file size limitations. (Actually, 1 GB is just the default segment size. The segment size can be adjusted using the configuration option `--with-segsize` when building Postgres Pro.) In principle, free space map and visibility map forks could require multiple segments as well, though this is unlikely to happen in practice.

A table that has columns with potentially large entries will have an associated *TOAST* table, which is used for out-of-line storage of field values that are too large to keep in the table rows proper. `pg_class.reltoastrelid` links from a table to its TOAST table, if any. See [Section 66.2](#) for more information.

The contents of tables and indexes are discussed further in [Section 66.6](#).

Tablespaces make the scenario more complicated. Each user-defined tablespace has a symbolic link inside the `PGDATA/pg_tblspc` directory, which points to the physical tablespace directory (i.e., the location specified in the tablespace's `CREATE TABLESPACE` command). This symbolic link is named after the tablespace's OID. Inside the physical tablespace directory there is a subdirectory with a name that depends on the Postgres Pro server version, such as `PG_9.0_201008051`. (The reason for using this subdirectory is so that successive versions of the database can use the same `CREATE TABLESPACE` location value without conflicts.) Within the version-specific subdirectory, there is a subdirectory for each database that has elements in the tablespace, named after the database's OID. Tables and indexes are stored within that directory, using the filenode naming scheme. The `pg_default` tablespace is not accessed through `pg_tblspc`, but corresponds to `PGDATA/base`. Similarly, the `pg_global` tablespace is not accessed through `pg_tblspc`, but corresponds to `PGDATA/global`.

The `pg_relation_filepath()` function shows the entire path (relative to `PGDATA`) of any relation. It is often useful as a substitute for remembering many of the above rules. But keep in mind that this function just gives the name of the first segment of the main fork of the relation — you may need to append a segment number and/or `_fsm`, `_vm`, or `_init` to find all the files associated with the relation.

Temporary files (for operations such as sorting more data than can fit in memory) are created within `PGDATA/base/pgsql_tmp`, or within a `pgsql_tmp` subdirectory of a tablespace directory if a tablespace other than `pg_default` is specified for them. The name of a temporary file has the form `pgsql_tmpPPP-P.NNN`, where `PPP` is the PID of the owning backend and `NNN` distinguishes different temporary files of that backend.

66.2. TOAST

This section provides an overview of TOAST (The Oversized-Attribute Storage Technique).

Postgres Pro uses a fixed page size (commonly 8 kB), and does not allow tuples to span multiple pages. Therefore, it is not possible to store very large field values directly. To overcome this limitation, large field values are compressed and/or broken up into multiple physical rows. This happens transparently to the user, with only small impact on most of the backend code. The technique is affectionately known as TOAST (or “the best thing since sliced bread”). The TOAST infrastructure is also used to improve handling of large data values in-memory.

Only certain data types support TOAST — there is no need to impose the overhead on data types that cannot produce large field values. To support TOAST, a data type must have a variable-length (*varlena*) representation, in which, ordinarily, the first four-byte word of any stored value contains the total length of the value in bytes (including itself). TOAST does not constrain the rest of the data type's representation. The special representations collectively called *TOASTed values* work by modifying or reinterpreting this initial length word. Therefore, the C-level functions supporting a TOAST-able data type must be careful about how they handle potentially TOASTed input values: an input might not actually consist of a four-byte length word and contents until after it's been *detoasted*. (This is normally done by invoking `PG_DETOAST_DATUM` before doing anything with an input value, but in some cases more efficient approaches are possible. See [Section 39.13.1](#) for more detail.)

TOAST usurps two bits of the *varlena* length word (the high-order bits on big-endian machines, the low-order bits on little-endian machines), thereby limiting the logical size of any value of a TOAST-able data type to 1 GB (2^{30} - 1 bytes). When both bits are zero, the value is an ordinary un-TOASTed value of the data type, and the remaining bits of the length word give the total datum size (including length word) in bytes. When the highest-order or lowest-order bit is set, the value has only a single-byte header instead of the normal four-byte header, and the remaining bits of that byte give the total datum size (including length byte) in bytes. This alternative supports space-efficient storage of values shorter than 127 bytes, while still allowing the data type to grow to 1 GB at need. Values with single-byte headers aren't aligned on any particular boundary, whereas values with four-byte headers are aligned on at least a four-byte boundary; this omission of alignment padding provides additional space savings that is significant compared to short values. As a special case, if the remaining bits of a single-byte header are all zero (which would be impossible for a self-inclusive length), the value is a pointer to out-of-line

data, with several possible alternatives as described below. The type and size of such a *TOAST pointer* are determined by a code stored in the second byte of the datum. Lastly, when the highest-order or lowest-order bit is clear but the adjacent bit is set, the content of the datum has been compressed and must be decompressed before use. In this case the remaining bits of the four-byte length word give the total size of the compressed datum, not the original data. Note that compression is also possible for out-of-line data but the varlena header does not tell whether it has occurred — the content of the TOAST pointer tells that, instead.

The compression technique used for either in-line or out-of-line compressed data can be selected for each column by setting the `COMPRESSION` column option in `CREATE TABLE` or `ALTER TABLE`. The default for columns with no explicit setting is to consult the `default_toast_compression` parameter at the time data is inserted.

As mentioned, there are multiple types of TOAST pointer datums. The oldest and most common type is a pointer to out-of-line data stored in a *TOAST table* that is separate from, but associated with, the table containing the TOAST pointer datum itself. These *on-disk* pointer datums are created by the TOAST management code when a tuple to be stored on disk is too large to be stored as-is. Further details appear in [Section 66.2.1](#). Alternatively, a TOAST pointer datum can contain a pointer to out-of-line data that appears elsewhere in memory. Such datums are necessarily short-lived, and will never appear on-disk, but they are very useful for avoiding copying and redundant processing of large data values. Further details appear in [Section 66.2.2](#).

66.2.1. Out-of-Line, On-Disk TOAST Storage

If any of the columns of a table are TOAST-able, the table will have an associated TOAST table, whose OID is stored in the table's `pg_class.reltoastrelid` entry. On-disk TOASTed values are kept in the TOAST table, as described in more detail below.

Out-of-line values are divided (after compression if used) into chunks of at most `TOAST_MAX_CHUNK_SIZE` bytes (by default this value is chosen so that four chunk rows will fit on a page, making it about 2000 bytes). Each chunk is stored as a separate row in the TOAST table belonging to the owning table. Every TOAST table has the columns `chunk_id` (an OID identifying the particular TOASTed value), `chunk_seq` (a sequence number for the chunk within its value), and `chunk_data` (the actual data of the chunk). A unique index on `chunk_id` and `chunk_seq` provides fast retrieval of the values. A pointer datum representing an out-of-line on-disk TOASTed value therefore needs to store the OID of the TOAST table in which to look and the OID of the specific value (its `chunk_id`). For convenience, pointer datums also store the logical datum size (original uncompressed data length), physical stored size (different if compression was applied), and the compression method used, if any. Allowing for the varlena header bytes, the total size of an on-disk TOAST pointer datum is therefore 18 bytes regardless of the actual size of the represented value.

The TOAST management code is triggered only when a row value to be stored in a table is wider than `TOAST_TUPLE_THRESHOLD` bytes (normally 2 kB). The TOAST code will compress and/or move field values out-of-line until the row value is shorter than `TOAST_TUPLE_TARGET` bytes (also normally 2 kB, adjustable) or no more gains can be had. During an `UPDATE` operation, values of unchanged fields are normally preserved as-is; so an `UPDATE` of a row with out-of-line values incurs no TOAST costs if none of the out-of-line values change.

The TOAST management code recognizes four different strategies for storing TOAST-able columns on disk:

- `PLAIN` prevents either compression or out-of-line storage. This is the only possible strategy for columns of non-TOAST-able data types.
- `EXTENDED` allows both compression and out-of-line storage. This is the default for most TOAST-able data types. Compression will be attempted first, then out-of-line storage if the row is still too big.
- `EXTERNAL` allows out-of-line storage but not compression. Use of `EXTERNAL` will make substring operations on wide `text` and `bytea` columns faster (at the penalty of increased storage space) because these operations are optimized to fetch only the required parts of the out-of-line value when it is not compressed.

- `MAIN` allows compression but not out-of-line storage. (Actually, out-of-line storage will still be performed for such columns, but only as a last resort when there is no other way to make the row small enough to fit on a page.)

Each TOASTable data type specifies a default strategy for columns of that data type, but the strategy for a given table column can be altered with `ALTER TABLE ... SET STORAGE`.

`TOAST_TUPLE_TARGET` can be adjusted for each table using `ALTER TABLE ... SET (toast_tuple_target = N)`

This scheme has a number of advantages compared to a more straightforward approach such as allowing row values to span pages. Assuming that queries are usually qualified by comparisons against relatively small key values, most of the work of the executor will be done using the main row entry. The big values of TOASTed attributes will only be pulled out (if selected at all) at the time the result set is sent to the client. Thus, the main table is much smaller and more of its rows fit in the shared buffer cache than would be the case without any out-of-line storage. Sort sets shrink also, and sorts will more often be done entirely in memory. A little test showed that a table containing typical HTML pages and their URLs was stored in about half of the raw data size including the TOAST table, and that the main table contained only about 10% of the entire data (the URLs and some small HTML pages). There was no run time difference compared to an un-TOASTed comparison table, in which all the HTML pages were cut down to 7 kB to fit.

66.2.2. Out-of-Line, In-Memory TOAST Storage

TOAST pointers can point to data that is not on disk, but is elsewhere in the memory of the current server process. Such pointers obviously cannot be long-lived, but they are nonetheless useful. There are currently two sub-cases: pointers to *indirect* data and pointers to *expanded* data.

Indirect TOAST pointers simply point at a non-indirect varlena value stored somewhere in memory. This case was originally created merely as a proof of concept, but it is currently used during logical decoding to avoid possibly having to create physical tuples exceeding 1 GB (as pulling all out-of-line field values into the tuple might do). The case is of limited use since the creator of the pointer datum is entirely responsible that the referenced data survives for as long as the pointer could exist, and there is no infrastructure to help with this.

Expanded TOAST pointers are useful for complex data types whose on-disk representation is not especially suited for computational purposes. As an example, the standard varlena representation of a PostgreSQL array includes dimensionality information, a nulls bitmap if there are any null elements, then the values of all the elements in order. When the element type itself is variable-length, the only way to find the *N*th element is to scan through all the preceding elements. This representation is appropriate for on-disk storage because of its compactness, but for computations with the array it's much nicer to have an "expanded" or "deconstructed" representation in which all the element starting locations have been identified. The TOAST pointer mechanism supports this need by allowing a pass-by-reference Datum to point to either a standard varlena value (the on-disk representation) or a TOAST pointer that points to an expanded representation somewhere in memory. The details of this expanded representation are up to the data type, though it must have a standard header and meet the other API requirements given in `src/include/utils/expandeddatum.h`. C-level functions working with the data type can choose to handle either representation. Functions that do not know about the expanded representation, but simply apply `PG_DETOAST_DATUM` to their inputs, will automatically receive the traditional varlena representation; so support for an expanded representation can be introduced incrementally, one function at a time.

TOAST pointers to expanded values are further broken down into *read-write* and *read-only* pointers. The pointed-to representation is the same either way, but a function that receives a read-write pointer is allowed to modify the referenced value in-place, whereas one that receives a read-only pointer must not; it must first create a copy if it wants to make a modified version of the value. This distinction and some associated conventions make it possible to avoid unnecessary copying of expanded values during query execution.

For all types of in-memory TOAST pointer, the TOAST management code ensures that no such pointer datum can accidentally get stored on disk. In-memory TOAST pointers are automatically expanded to

normal in-line varlena values before storage — and then possibly converted to on-disk TOAST pointers, if the containing tuple would otherwise be too big.

66.2.3. Pluggable TOAST

TOAST is a part of Postgres Pro core. TOAST is not extensible, has the same strategy for all data types and is not effective for structured data (JSON) or data that requires a special workflow (`bytea`).

Experimental pluggable TOAST provides an open API allowing you to develop and plug in custom TOAST implementations for table columns and data types in addition to the default one. These implementations are called TOASTers. A TOASTER can process TOASTed data using the knowledge about the internal data structure and workflow.

TOASTers can use any storage, advanced compression, encryption and other functionalities internally, without being aware of the table access method that calls the TOASTER.

The pluggable TOAST consists of the core, that is, the TOAST API, and custom TOASTers plugged in to Postgres Pro with this API.

This API does not possess any logic regarding the TOAST functionality, it is just a wraparound that allows registering, calling and dropping plugged TOASTers.

To enable using custom TOAST implementations, the TOAST API defines a new `Custom` type of TOAST pointer in addition to `External` and `Extended` (corresponding to `EXTERNAL` and `EXTENDED` TOAST-able column storing strategies).

Currently the pluggable TOAST has several limitations:

- Custom TOASTers must only be assigned to a column with the `EXTERNAL` storing strategy.
- Currently the pluggable TOAST does not support logical replication. A custom TOAST pointer is skipped by the logical replication engine and is not replicated. This is because support of custom pointers requires heavy modifications in the Postgres Pro core (the logical replication engine) and even architecture changes as now custom pointers allow updating TOASTed values, while default pointers do not.
- Custom TOASTers are not restored from dumps and require manual actions for restoring. The TOASTER extension is not restored from dump, thus the data processed with it are not available. Another point is that the same TOASTER extension has different IDs in different systems, so restoring a TOASTed value from dump would require determining the TOASTER used for processing data.
- When a custom TOASTER is assigned to a table column, in all the transactions after the TOASTER assignment, data being stored in the TOAST tables is processed with the TOASTER assigned to the column.
- When a custom TOASTER is assigned to a table column, old data values are left as-is and are not re-TOASTed although while being updated, an old data value is TOASTed with the assigned TOASTER because of the workflow of the default TOAST mechanics.
- If the TOASTER is somehow dropped, the data TOASTed with it will not be accessible with any standard techniques and is considered to be lost. To protect TOASTed values, the API creates some dependencies in the system (`PG_DEPEND` table):
 - When the TOASTER is added to a system, a TOASTER-extension dependency is created.
 - When the TOASTER is assigned to a column, a TOASTER-relation dependency is created.
 - When the TOASTER is unassigned from a column, the TOASTER-relation dependency is deleted.
 - When the TOASTER is dropped from a system, the TOASTER-extension dependency is deleted.

Custom TOASTers are added as extensions implementing some mandatory functions (check `README.toastapi` file in `/contrib/toastapi` for more details). Pluggable TOAST provides a set of SQL functions to control custom TOASTers. These functions are called with SQL `SELECT` command:

`add_toaster(toaster_name text, toaster_handler_func text) → integer`

Plug in a new TOASTer (added before as an extension). This function also creates a dependency in `PG_DEPEND` for the TOASTer OID to its extension OID, which is used to protect the TOASTer from being dropped when dropping the TOASTer extension.

Arguments:

- `toaster_name` — TOASTer name to be stored in `PG_TOASTER`. Must be unique.
- `toaster_handler_func` — TOASTer handler function name (automatically added as a `regproc` when the TOASTer extension is installed).

Returns OID assigned to the new TOASTer.

`set_toaster(toaster_name text, tab_name text, col_name text) → integer`

Assign a TOASTer to a table column. This function also creates a dependency in `PG_DEPEND` for the TOASTer OID to the relation OID. This dependency is used to protect TOASTed data by restricting dropping the TOASTer that is assigned to a column.

Arguments:

- `toaster_name` — name of the TOASTer assigned to a table column.
- `tab_name` — table name.
- `col_name` — column name.

Returns OID of the TOASTer assigned to the column.

`reset_toaster(tab_name text, col_name text) → integer`

Reset the TOAST mechanism to a default one for a table column. This function drops the TOASTer-relation dependency, which allows dropping the TOASTer later, when all dependencies are dropped.

Arguments:

- `tab_name` — table name.
- `col_name` — table column name.

Returns always 0.

`drop_toaster(toaster_name text)`

Drop a TOASTer from the system. Only unused TOASTers, which do not have any data TOASTed with them can be dropped. Dropping is protected by the dependencies created by `add_toaster` and `set_toaster` functions. If a dependency still exists, the TOASTer would not be dropped. If the TOASTer was dropped successfully, the extension with its implementation can be dropped by the `DROP EXTENSION` command.

Arguments:

- `toaster_name` — name of the TOASTer to be dropped.

Returns the TOASTer OID on success or 0 if the TOASTer is used.

`get_toaster(tab_name text, col_name text) → integer`

Get the OID of the TOASTer assigned to a table column.

Arguments:

- `tab_name` — table name.
- `col_name` — table column name.

Returns the OID of the TOASTer assigned to a column or 0 if none.

The following example illustrates the use of the pluggable TOAST API:

```
CREATE EXTENSION toastapi;

CREATE EXTENSION bytea_toaster;

SELECT add_toaster('bytea_toaster');

CREATE TABLE test_bytea_append (id int, a bytea);

ALTER TABLE test_bytea_append ALTER a SET STORAGE external;

SELECT set_toaster('bytea_toaster', 'test_bytea_append', 'a');
       test_set_toaster
-----
(1 row)

...

SELECT get_toaster('test_bytea_append', 'a') AS bytea_toaster_oid;
       get_toaster
-----
          16348
(1 row)

SELECT pgpro_toast.reset_toaster('test_bytea_append', 'a');
       reset_toaster
-----
              0
(1 row)

DROP TABLE test_bytea_append;
SELECT pgpro_toast.drop_toaster('bytea_toaster');
       drop_toaster
-----
          16348
(1 row)

DROP EXTENSION bytea_toaster;
```

66.3. Free Space Map

Each heap and index relation, except for hash indexes, has a Free Space Map (FSM) to keep track of available space in the relation. It's stored alongside the main relation data in a separate relation fork, named after the filenode number of the relation, plus a `_fsm` suffix. For example, if the filenode of a relation is 12345, the FSM is stored in a file called `12345_fsm`, in the same directory as the main relation file.

The Free Space Map is organized as a tree of FSM pages. The bottom level FSM pages store the free space available on each heap (or index) page, using one byte to represent each such page. The upper levels aggregate information from the lower levels.

Within each FSM page is a binary tree, stored in an array with one byte per node. Each leaf node represents a heap page, or a lower level FSM page. In each non-leaf node, the higher of its children's values is stored. The maximum value in the leaf nodes is therefore stored at the root.

66.4. Visibility Map

Each heap relation has a Visibility Map (VM) to keep track of which pages contain only tuples that are known to be visible to all active transactions; it also keeps track of which pages contain only frozen tuples. It's stored alongside the main relation data in a separate relation fork, named after the filenode number of the relation, plus a `_vm` suffix. For example, if the filenode of a relation is 12345, the VM is stored in a file called `12345_vm`, in the same directory as the main relation file. Note that indexes do not have VMs.

The visibility map stores two bits per heap page. The first bit, if set, indicates that the page is all-visible, or in other words that the page does not contain any tuples that need to be vacuumed. This information can also be used by *index-only scans* to answer queries using only the index tuple. The second bit, if set, means that all tuples on the page have been frozen.

The map is conservative in the sense that we make sure that whenever a bit is set, we know the condition is true, but if a bit is not set, it might or might not be true. Visibility map bits are only set by vacuum, but are cleared by any data-modifying operations on a page.

66.5. The Initialization Fork

Each unlogged table, and each index on an unlogged table, has an initialization fork. The initialization fork is an empty table or index of the appropriate type. When an unlogged table must be reset to empty due to a crash, the initialization fork is copied over the main fork, and any other forks are erased (they will be recreated automatically as needed).

66.6. Database Page Layout

This section provides an overview of the page format used within Postgres Pro tables and indexes.¹ Sequences and TOAST tables are formatted just like a regular table.

In the following explanation, a *byte* is assumed to contain 8 bits. In addition, the term *item* refers to an individual data value that is stored on a page. In a table, an item is a row; in an index, an item is an index entry.

Every table and index is stored as an array of *pages* of a fixed size (usually 8 kB, although a different page size can be selected when compiling the server). In a table, all the pages are logically equivalent, so a particular item (row) can be stored in any page. In indexes, the first page is generally reserved as a *metapage* holding control information, and there can be different types of pages within the index, depending on the index access method.

Table 66.2 shows the overall layout of a page. There are five parts to each page.

Table 66.2. Overall Page Layout

Item	Description
PageHeaderData	24 bytes long. Contains general information about the page, including free space pointers.
ItemIdData	Array of item identifiers pointing to the actual items. Each entry is an (offset,length) pair. 4 bytes per item.
Free space	The unallocated space. New item identifiers are allocated from the start of this area, new items from the end.
Items	The actual items themselves.

¹ Actually, use of this page format is not required for either table or index access methods. The `heap` table access method always uses this format. All the existing index methods also use the basic format, but the data kept on index metapages usually doesn't follow the item layout rules.

Item	Description
Special space	Index access method specific data. Different methods store different data. Additionally, this space is used to store information about 64-bit transaction IDs.

The first 20 bytes of each page consists of a page header (`PageHeaderData`). Its format is detailed in [Table 66.3](#). The first field tracks the most recent WAL entry related to this page. The second field contains the page checksum if [data checksums](#) are enabled. Next is a 2-byte field containing flag bits. This is followed by three 2-byte integer fields (`pd_lower`, `pd_upper`, and `pd_special`). These contain byte offsets from the page start to the start of unallocated space, to the end of unallocated space, and to the start of the special space. The next 2 bytes of the page header, `pd_pagesize_version`, store both the page size and a version indicator. Beginning with PostgreSQL 8.3 the version number is 4; PostgreSQL 8.1 and 8.2 used version number 3; PostgreSQL 8.0 used version number 2; PostgreSQL 7.3 and 7.4 used version number 1; prior releases used version number 0. (The basic page layout and header format has not changed in most of these versions, but the layout of heap row headers has.) The page size is basically only present as a cross-check; there is no support for having more than one page size in an installation.

Table 66.3. PageHeaderData Layout

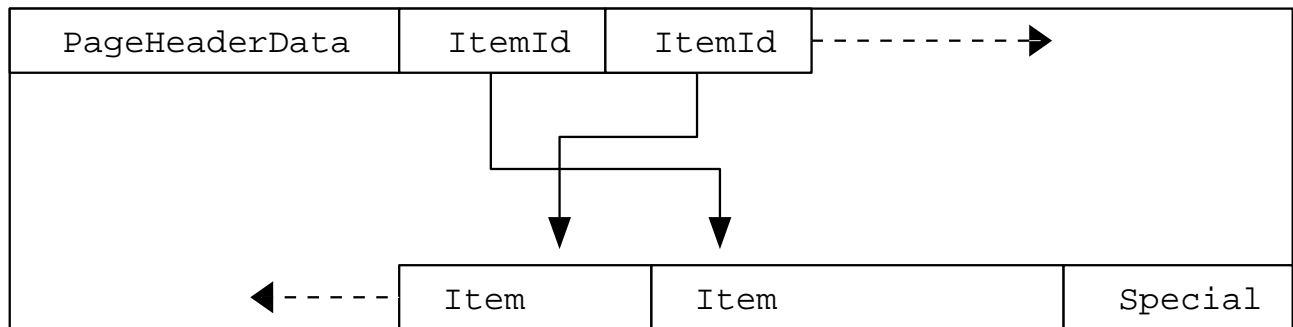
Field	Type	Length	Description
<code>pd_lsn</code>	<code>PageXLogRecPtr</code>	8 bytes	LSN: next byte after last byte of WAL record for last change to this page
<code>pd_checksum</code>	<code>uint16</code>	2 bytes	Page checksum
<code>pd_flags</code>	<code>uint16</code>	2 bytes	Flag bits
<code>pd_lower</code>	<code>LocationIndex</code>	2 bytes	Offset to start of free space
<code>pd_upper</code>	<code>LocationIndex</code>	2 bytes	Offset to end of free space
<code>pd_special</code>	<code>LocationIndex</code>	2 bytes	Offset to start of special space
<code>pd_pagesize_version</code>	<code>uint16</code>	2 bytes	Page size and layout version number information

Following the page header are item identifiers (`ItemIdData`), each requiring four bytes. An item identifier contains a byte-offset to the start of an item, its length in bytes, and a few attribute bits which affect its interpretation. New item identifiers are allocated as needed from the beginning of the unallocated space. The number of item identifiers present can be determined by looking at `pd_lower`, which is increased to allocate a new identifier. Because an item identifier is never moved until it is freed, its index can be used on a long-term basis to reference an item, even when the item itself is moved around on the page to compact free space. In fact, every pointer to an item (`ItemPointer`, also known as `CTID`) created by Postgres Pro consists of a page number and the index of an item identifier.

The items themselves are stored in space allocated backwards from the end of unallocated space. The exact structure varies depending on what the table is to contain. Tables and sequences both use a structure named `HeapTupleHeaderData`, described below.

The final section is the “special section” which can contain anything the access method wishes to store. For example, b-tree indexes store links to the page's left and right siblings, as well as some other data relevant to the index structure. Ordinary tables and sequences store the `HeapPageSpecialData` structure in this section. Its format is specified in [Table 66.4](#).

[Figure 66.1](#) illustrates how these parts are laid out in a page.

Figure 66.1. Page Layout

66.6.1. Table Row Layout

`HeapPageSpecialData` is a 24-byte structure. The first two 8-byte fields contain the base (i.e., an offset) for short (4-byte) transaction IDs in this page. The next field is a hint that shows whether pruning the page is likely to be profitable: it tracks the oldest un-pruned XMAX on the page. The last field is a magic number identifying the type of the page (0x1010 for ordinary tables and 0x1717 for sequences).

Table 66.4. HeapPageSpecialData Layout

Field	Type	Length	Description
pd_xid_base	TransactionId	8 bytes	Base for short 4-byte transaction IDs in this page
pd_multi_base	TransactionId	8 bytes	Base for short 4-byte multixact IDs in this page
pd_prune_xid	ShortTransactionId	4 bytes	Oldest unpruned XMAX on page, or zero if none
pd_magic	uint32	4 bytes	Magic number identifying the type of this page

All table rows are structured in the same way. There is a fixed-size header (occupying 40 bytes on most machines), followed by an optional null bitmap, an optional object ID field, and the user data. The header is detailed in [Table 66.5](#). The actual user data (columns of the row) begins at the offset indicated by `t_hoff`, which must always be a multiple of the MAXALIGN distance for the platform. The null bitmap is only present if the `HEAP_HASNULL` bit is set in `t_infomask`. If it is present it begins just after the fixed header and occupies enough bytes to have one bit per data column (that is, the number of bits that equals the attribute count in `t_infomask2`). In this list of bits, a 1 bit indicates not-null, a 0 bit is a null. When the bitmap is not present, all columns are assumed not-null. The object ID is only present if the `HEAP_HASOID_OLD` bit is set in `t_infomask`. If present, it appears just before the `t_hoff` boundary. Any padding needed to make `t_hoff` a MAXALIGN multiple will appear between the null bitmap and the object ID. (This in turn ensures that the object ID is suitably aligned.)

Table 66.5. HeapTupleHeaderData Layout

Field	Type	Length	Description
t_xmin	ShortTransactionId	4 bytes	insert XID stamp
t_xmax	ShortTransactionId	4 bytes	delete XID stamp

Field	Type	Length	Description
t_cid	CommandId	4 bytes	insert and/or delete CID stamp (overlays with t_xvac)
t_xvac	ShortTransactionId	4 bytes	XID for VACUUM operation moving a row version
t_ctid	ItemPointerData	6 bytes	current TID of this or newer row version
t_infomask2	uint16	2 bytes	number of attributes, plus various flag bits
t_infomask	uint16	2 bytes	various flag bits
t_hoff	uint8	1 byte	offset to user data

Interpreting the actual data can only be done with information obtained from other tables, mostly `pg_attribute`. The key values needed to identify field locations are `attlen` and `attalign`. There is no way to directly get a particular attribute, except when there are only fixed width fields and no null values. All this trickery is wrapped up in the functions `heap_getattr`, `fastgetattr` and `heap_getsysattr`.

To read the data you need to examine each attribute in turn. First check whether the field is NULL according to the null bitmap. If it is, go to the next. Then make sure you have the right alignment. If the field is a fixed width field, then all the bytes are simply placed. If it's a variable length field (`attlen = -1`) then it's a bit more complicated. All variable-length data types share the common header structure `struct varlena`, which includes the total length of the stored value and some flag bits. Depending on the flags, the data can be either inline or in a TOAST table; it might be compressed, too (see [Section 66.2](#)).

66.7. Heap-Only Tuples (HOT)

To allow for high concurrency, Postgres Pro uses [multiversion concurrency control](#) (MVCC) to store rows. However, MVCC has some downsides for update queries. Specifically, updates require new versions of rows to be added to tables. This can also require new index entries for each updated row, and removal of old versions of rows and their index entries can be expensive.

To help reduce the overhead of updates, Postgres Pro has an optimization called heap-only tuples (HOT). This optimization is possible when:

- The update does not modify any columns referenced by the table's indexes, not including summarizing indexes. The only summarizing index method in the core Postgres Pro distribution is [BRIN](#).
- There is sufficient free space on the page containing the old row for the updated row.

In such cases, heap-only tuples provide two optimizations:

- New index entries are not needed to represent updated rows, however, summary indexes may still need to be updated.
- When a row is updated multiple times, row versions other than the oldest and the newest can be completely removed during normal operation, including `SELECTs`, instead of requiring periodic vacuum operations. (Indexes always refer to the [page item identifier](#) of the original row version. The tuple data associated with that row version is removed, and its item identifier is converted to a redirect that points to the oldest version that may still be visible to some concurrent transaction. Intermediate row versions that are no longer visible to anyone are completely removed, and the associated page item identifiers are made available for reuse.)

You can increase the likelihood of sufficient page space for HOT updates by decreasing a table's [fill-factor](#). If you don't, HOT updates will still happen because new rows will naturally migrate to new pages and existing pages with sufficient free space for new row versions. The system view [pg_stat_all_tables](#) allows monitoring of the occurrence of HOT and non-HOT updates.

Chapter 67. Transaction Processing

This chapter provides an overview of the internals of Postgres Pro's transaction management system. The word transaction is often abbreviated as *xact*.

67.1. Transactions and Identifiers

Transactions can be created explicitly using `BEGIN` or `START TRANSACTION` and ended using `COMMIT` or `ROLLBACK`. SQL statements outside of explicit transactions automatically use single-statement transactions.

Every transaction is identified by a unique `VirtualTransactionId` (also called `virtualXID` or `vxid`), which is comprised of a backend's process number (or `procNumber`) and a sequentially-assigned number local to each backend, known as `localXID`. For example, the virtual transaction ID 4/12532 has a `procNumber` of 4 and a `localXID` of 12532.

Non-virtual `TransactionIds` (or `xid`), e.g., 278394, are assigned sequentially to transactions from a global counter used by all databases within the Postgres Pro cluster. This assignment happens when a transaction first writes to the database. This means lower-numbered `xids` started writing before higher-numbered `xids`. Note that the order in which transactions perform their first database write might be different from the order in which the transactions started, particularly if the transaction started with statements that only performed database reads.

In Postgres Pro Shardman, the internal transaction ID type `xid` is 64 bits wide to prevent [transaction ID wraparound](#). Each tuple header contains two XIDs, so extending them would lead to high overhead. For that reason, when saved on disk on-tuple XIDs are 32-bit, but each page special area contains an offset, called *base XID*. When a tuple is read into memory, the base XID is added to its 32-bit XIDs, and both XIDs in the tuple become full 64 bits wide so the tuple becomes the so-called in-memory tuple. Full 64-bit XIDs acquired after adding the base XID are used in comparison and other arithmetic operations. In some contexts, `xid8` is also used. The functions in [Table 9.82](#) return `xid8` values. Xids are used as the basis for Postgres Pro's [MVCC](#) concurrency mechanism and streaming replication.

When a top-level transaction with a (non-virtual) `xid` commits, it is marked as committed in the `pg_xact` directory. Additional information is recorded in the `pg_commit_ts` directory if [track_commit_timestamp](#) is enabled.

In addition to `vxid` and `xid`, prepared transactions are also assigned Global Transaction Identifiers (GID). GIDs are string literals up to 200 bytes long, which must be unique amongst other currently prepared transactions. The mapping of GID to `xid` is shown in [pg_prepared_xacts](#).

67.2. Transactions and Locking

The transaction IDs of currently executing transactions are shown in [pg_locks](#) in columns `virtualxid` and `transactionid`. Read-only transactions will have `virtualxids` but NULL `transactionids`, while both columns will be set in read-write transactions.

Some lock types wait on `virtualxid`, while other types wait on `transactionid`. Row-level read locks might also require the assignment of multixact IDs (`mxid`; see [Section 24.1.5.1](#)).

67.3. Subtransactions

Subtransactions are started inside transactions, allowing large transactions to be broken into smaller units. Subtransactions can commit or abort without affecting their parent transactions, allowing parent transactions to continue. This allows errors to be handled more easily, which is a common application development pattern. The word subtransaction is often abbreviated as *subxact*.

Subtransactions can be started explicitly using the `SAVEPOINT` command, but can also be started in other ways, such as PL/pgSQL's `EXCEPTION` clause. PL/Python and PL/Tcl also support explicit subtransactions.

Subtransactions can also be started from other subtransactions. The top-level transaction and its child subtransactions form a hierarchy or tree, which is why we refer to the main transaction as the top-level transaction.

If a subtransaction is assigned a non-virtual transaction ID, its transaction ID is referred to as a “subxid”. In Postgres Pro Shardman, subxids are 64 bits wide. Read-only subtransactions are not assigned subxids, but once they attempt to write, they will be assigned one. This also causes all of a subxid’s parents, up to and including the top-level transaction, to be assigned non-virtual transaction ids. We ensure that a parent xid is always lower than any of its child subxids.

The immediate parent xid of each subxid is recorded in the `pg_subtrans` directory. No entry is made for top-level xids since they do not have a parent, nor is an entry made for read-only subtransactions.

When a subtransaction commits, all of its committed child subtransactions with subxids will also be considered subcommitted in that transaction. When a subtransaction aborts, all of its child subtransactions will also be considered aborted.

When a top-level transaction with an xid commits, all of its subcommitted child subtransactions are also persistently recorded as committed in the `pg_xact` subdirectory. If the top-level transaction aborts, all its subtransactions are also aborted, even if they were subcommitted.

The more subtransactions each transaction keeps open (not rolled back or released), the greater the transaction management overhead. Up to 64 open subxids are cached in shared memory for each back-end; after that point, the storage I/O overhead increases significantly due to additional lookups of subxid entries in `pg_subtrans`.

67.4. Two-Phase Transactions

Postgres Pro supports a two-phase commit (2PC) protocol that allows multiple distributed systems to work together in a transactional manner. The commands are `PREPARE TRANSACTION`, `COMMIT PREPARED` and `ROLLBACK PREPARED`. Two-phase transactions are intended for use by external transaction management systems. Postgres Pro follows the features and model proposed by the X/Open XA standard, but does not implement some less often used aspects.

When the user executes `PREPARE TRANSACTION`, the only possible next commands are `COMMIT PREPARED` or `ROLLBACK PREPARED`. In general, this prepared state is intended to be of very short duration, but external availability issues might mean transactions stay in this state for an extended interval. Short-lived prepared transactions are stored only in shared memory and WAL. Transactions that span checkpoints are recorded in the `pg_twophase` directory. Transactions that are currently prepared can be inspected using `pg_prepared_xacts`.

67.5. Distributed Transactions

67.5.1. Visibility and CSN

67.5.1.1. CSN — Commit Sequence Number

A Postgres Pro Shardman cluster uses a snapshot isolation mechanism for distributed transactions. The mechanism provides a way to synchronize snapshots between different nodes of a cluster and a way to atomically commit such a transaction with respect to other concurrent global and local transactions. These global transactions can be coordinated by using provided SQL functions or through `postgres_fdw`, which uses these functions on remote nodes transparently.

Assume that each node uses the CSN-based visibility: the database tracks the counter for each transaction commit (CSN). With such a setting, a snapshot is just a single number — a copy of the current CSN at the moment when the snapshot was taken. Visibility rules are boiled down to checking whether the current tuple’s CSN is less than our snapshot’s CSN.

Let's assume that `CSN` is the current physical time on the node and call it `GlobalCSN`. If the physical time on different nodes is perfectly synchronized, then such a snapshot obtained on one node can be used on other nodes to provide the necessary level of transaction isolation. But unfortunately physical time is never perfectly sync and can drift, and this should be taken into account. Also, there is no easy notion of lock or atomic operation in the distributed environment, so commit atomicity on different nodes with respect to concurrent snapshot acquisition should be handled somehow. This is addressed in the following way:

1. To achieve commit atomicity of different nodes, intermediate step is introduced: at the first run, a transaction is marked as `InDoubt` on all nodes, and only after that each node commits it and stamps with a given `GlobalCSN`. All readers that ran into tuples of an `InDoubt` transaction should wait until it ends and recheck the visibility.
2. When the coordinator is marking transactions as `InDoubt` on other nodes, it collects `ProposedGlobalCSN` from each participant, which is the local time on those nodes. Next, it selects the maximal value of all `ProposedGlobalCSNs` and commits the transaction on all nodes with that maximal `GlobalCSN` even if that value is greater than the current time on this node due to clock drift. So the `GlobalCSN` for the given transaction will be the same on all nodes. Each node records its last generated `CSN` (`last_csn`) and cannot generate `CSN ≤ last_csn`. When a node commits a transaction with `CSN > last_csn`, `last_csn` is adjusted to record this `CSN`. Due to this mechanism, a node cannot generate a `CSN`, that is less than `CSNs` of already committed transactions.
3. When a local transaction imports a foreign global snapshot with some `GlobalCSN` and the current time on this node is smaller than the incoming `GlobalCSN`, then the transaction should wait until this `GlobalCSN` time comes to the local clock.

The two last rules provide protection against time drift.

67.5.1.2. Commit Delay and External Consistency

The rules above still do not guarantee recency for snapshots generated on nodes that do not participate in a transaction. A read operation that originates from such a node can see stale data. The probability of the anomaly directly depends on the system clock skew in the Postgres Pro Shardman cluster.

Particular attention should be paid to the synchronization of system clocks on all cluster nodes. The size of the clock skew must be measured. If an external consistency is required, then the clock skew can be compensated with a commit delay. This delay is added before every commit in the system, so it has a negative impact on the latency of transactions. Read-only transactions are not affected by this delay. The delay can be set using the configuration parameter `csn_commit_delay`.

67.5.1.3. CSN Map

The `CSN` visibility mechanism described above is not a general way to check the visibility of all transactions. It is used to provide isolation only for distributed transactions. As a result, each cluster node uses a visibility checking mechanism based on `xid` and `xmin`. To be able to use the `CSN` snapshot that points to the past, we need to keep old versions of tuples on all nodes and therefore defer vacuuming them. To do this, each node in a Postgres Pro Shardman cluster maintains a `CSN` to `xid` mapping. The map is called `CSNSnapshotXidMap`. This map is a ring buffer, and it stores the correspondence between the current `snapshot_csn` and `xmin` in a sparse way: `snapshot_csn` is rounded to seconds (and here we use the fact that `snapshot_csn` is just a timestamp), and `xmin` is stored in the circular buffer where rounded `snapshot_csn` acts as an offset from the current circular buffer head. The size of the circular buffer is controlled by the `csn_snapshot_defer_time` configuration setting. `VACUUM` is not allowed to clean up tuples whose `xmax` is newer than the oldest `xmin` in `CSNSnapshotXidMap`.

When a `CSN` snapshot arrives, we check that its `snapshot_csn` is still in our map, otherwise, we will error out with "snapshot too old" message. If the `snapshot_csn` is successfully mapped, we fill backend's `xmin` with the value from the map. That way we can take into account backends with an imported `CSN` snapshot, and old tuple versions will be preserved.

67.5.1.4. CSN Map Trimming

To support global transactions, each node keeps old versions of tuples for at least `csn_snapshot_defer_time` seconds. With large values of `csn_snapshot_defer_time`, this negatively affects performance. This is because nodes save all row versions during the last `csn_snapshot_defer_time` seconds, but there may not be more transactions in the cluster that can read them. A special task of the monitor periodically recalculates `xmin` in the cluster and sets it on all nodes to the minimum possible value. This allows the vacuuming routine to remove a row version that is no longer of interest to any transaction. The `shardman.monitor.trim_csnxid_map_interval` configuration setting controls the worker. The worker wakes up every `monitor_interval` seconds and performs the following operations:

1. Checks if the current node's repgroup ID is the smallest among all IDs in the cluster. If this condition is not met, then the work on the current node is terminated. So only one node in the cluster can perform a horizon negotiation.
2. From each node of the Postgres Pro Shardman cluster, the coordinator collects the oldest snapshot CSN among all active transactions on the node.
3. The coordinator chooses the smallest CSN and sends it to each node. Each node discards its `csnXidMap` values that are less than this value.

67.5.2. 2PC and Prepared Transaction Resolution

Postgres Pro Shardman implements a two-phase commit protocol to ensure the atomicity of distributed transactions. During the execution of a distributed transaction, the coordinator node sends the command `BEGIN` to participant nodes to initiate their local transactions.

The term "participant nodes" herein and subsequently refers to a subset of cluster nodes that participate in the execution of a transaction's command while the node is engaged in writing activity.

Additionally, a local transaction is created on the coordinator node. This ensures that there are corresponding local transactions on all nodes participating in the distributed transaction.

During the two-phase transaction commit, the coordinator node sends the command `PREPARE TRANSACTION` to the participant nodes to initiate the preparation of their local transactions for commit. If the preparation is successful, the local transaction data is stored in a disk storage, making it persistent. If all participant nodes report successful preparation to the coordinator node, the coordinator node will commit its local transaction. Subsequently, the coordinator node will also commit the previously prepared transactions on the participant nodes using the command `COMMIT PREPARED`.

If a failure occurs during the `PREPARE TRANSACTION` command on any of the participant nodes, the distributed transaction is considered aborted. The coordinator node then broadcasts the command to abort the previously prepared transactions using the `ROLLBACK PREPARED` command. If the local transaction was already prepared, it is aborted. However, if there was no prepared transaction with the specified name, the command to rollback is simply ignored. Subsequently, the coordinator node rolls back its local transaction.

After a successful preparation phase, there will be an object `prepared transaction` on the each of participant nodes. These objects are actually disk files and records in the server memory.

It is possible to have a prepared transaction that was created earlier through a two-phase operation and will never be completed. This can occur, for example, if the coordinator node fails exactly after the preparation step but before the commit step. It can also occur as a result of network connectivity issues. For instance, if the command `COMMIT PREPARED` from the coordinator node to a participant node ends with an error, local transactions will be committed on all participant nodes except for the one with the error. The local transaction will also be committed on the coordinator node. All participants, except for the one with the error, believe that the distributed transaction was completed. However, the one participant still waiting for `COMMIT PREPARED` will never receive it, resulting in a prepared transaction that will never be completed.

A prepared transaction consumes system resources, such as memory and disk space. An incomplete prepared transaction causes other transactions that access rows modified by that transaction to wait until the distributed operation completes. Therefore, it is necessary to complete prepared transactions, even in cases where there were failures during commit, to free up resources and ensure that other transactions can proceed.

To resolve such situations, there is a mechanism for resolving prepared transactions that is implemented as part of the Postgres Pro Shardman monitor. It is implemented as a background worker that wakes up periodically, acting as an internal “crontab” job. By default, the period is set to 5 seconds, but it can be configured using the `shardman.monitor_dxact_interval` configuration parameter. The worker checks the presence of prepared transactions that were created earlier by a certain amount of time, specified by the `shardman.monitor_dxact_timeout` configuration parameter (which is also set to 5 seconds by default), on the same node where the Postgres Pro Shardman monitor is running.

When the `PREPARE TRANSACTION` command is sent to a participant node, a special name is assigned to the prepared transaction. This name encodes useful information, which allows identifying the coordinator node and its local transaction.

If the Postgres Pro Shardman monitor finds outdated prepared transactions, it extracts the coordinator's replication group ID and transaction ID of the coordinator's local transaction. The monitor then sends a query to the coordinator

```
SELECT shardman.xact_status(TransactionId)
```

which requests the current status of the coordinator's local transaction. If the query fails, for example, due to network connectivity issues, then the prepared transaction will remain untouched until the next time when the monitor wakes up.

In the case of a successful query, the coordinator node can reply with one of the following statuses:

`committed`

The local transaction on the coordinator node was completed successfully. Therefore, the Postgres Pro Shardman monitor also commits this prepared transaction using the `COMMIT PREPARED` command.

`aborted`

The local transaction on the coordinator node was aborted. Therefore, the monitor also aborts this transaction using the `ROLLBACK PREPARED` command.

`unknown`

The transaction with such an identifier never existed on the coordinator node. Therefore, the monitor aborts this transaction using the `ROLLBACK PREPARED` command.

`active`

The local transaction on the coordinator node is still somewhere inside the `CommitTransaction()` flow. Therefore, the monitor does nothing with this transaction. The monitor will try again with this transaction at the next wake-up.

`ambiguous`

This status can be returned when `CLOG`'s truncating is enabled on the coordinator node. The `CLOG` is a bitmap that stores the status of completed local transactions. When a transaction is committed or aborted, its status is marked in the `CLOG`. However, the `CLOG` can be truncated (garbage collected) by the `VACUUM` process to discard statuses of old transactions that do not affect the visibility of data for any existing transaction.

When the `CLOG` is truncated, there is a possibility that the `shardman.xact_status()` function may not be able to unambiguously decide if a transaction exists in the past (with some status) or if it never existed. In such cases, the function returns an ambiguous status. This can lead to uncertainty about the actual status of the transaction and can make it difficult to resolve the prepared transaction.

When the `shardman.xact_status()` function returns the `ambiguous` status for a prepared transaction, the monitor node logs a warning message indicating that the status could not be determined unambiguously. The prepared transaction is left untouched, and the monitor will try again with this transaction at the next wake-up. It is important to properly configure the `min_clog_size` parameter with the value of `1024000` (which means "never truncate CLOG") to avoid ambiguity in the status of prepared transactions.

In situations where the prepared transaction resolution mechanism is unable to resolve prepared transactions due to constant errors or ambiguous status, the administrator will need to manually intervene to resolve these transactions. This may involve examining the server logs and performing a manual rollback or commit operation on the prepared transaction. Note that leaving prepared transactions unresolved can lead to resource-consumption and performance issues, so it is important to address these situations as soon as possible.

Chapter 68. How the Planner Uses Statistics

This chapter builds on the material covered in [Section 14.1](#) and [Section 14.2](#) to show some additional details about how the planner uses the system statistics to estimate the number of rows each part of a query might return. This is a significant part of the planning process, providing much of the raw material for cost calculation.

The intent of this chapter is not to document the code in detail, but to present an overview of how it works. This will perhaps ease the learning curve for someone who subsequently wishes to read the code.

68.1. Row Estimation Examples

The examples shown below use tables in the Postgres Pro regression test database. Note also that since `ANALYZE` uses random sampling while producing statistics, the results will change slightly after any new `ANALYZE`.

Let's start with a very simple query:

```
EXPLAIN SELECT * FROM tenk1;
```

QUERY PLAN

```
-----  
Seq Scan on tenk1  (cost=0.00..458.00 rows=10000 width=244)
```

How the planner determines the cardinality of `tenk1` is covered in [Section 14.2](#), but is repeated here for completeness. The number of pages and rows is looked up in `pg_class`:

```
SELECT relpages, reltuples FROM pg_class WHERE relname = 'tenk1';
```

```
relpages | reltuples  
-----+-----
```

```
358 | 10000
```

These numbers are current as of the last `VACUUM` or `ANALYZE` on the table. The planner then fetches the actual current number of pages in the table (this is a cheap operation, not requiring a table scan). If that is different from `relpages` then `reltuples` is scaled accordingly to arrive at a current number-of-rows estimate. In the example above, the value of `relpages` is up-to-date so the rows estimate is the same as `reltuples`.

Let's move on to an example with a range condition in its `WHERE` clause:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 1000;
```

QUERY PLAN

```
-----  
Bitmap Heap Scan on tenk1  (cost=24.06..394.64 rows=1007 width=244)
```

```
Recheck Cond: (unique1 < 1000)
```

```
-> Bitmap Index Scan on tenk1_unique1  (cost=0.00..23.80 rows=1007 width=0)
```

```
Index Cond: (unique1 < 1000)
```

The planner examines the `WHERE` clause condition and looks up the selectivity function for the operator `<` in `pg_operator`. This is held in the column `oprrest`, and the entry in this case is `scalarlttsel`. The `scalarlttsel` function retrieves the histogram for `unique1` from `pg_statistic`. For manual queries it is more convenient to look in the simpler `pg_stats` view:

```
SELECT histogram_bounds FROM pg_stats
```

```
WHERE tablename='tenk1' AND attname='unique1';
```

```
histogram_bounds  
-----
```

```
{0, 993, 1997, 3050, 4040, 5036, 5957, 7057, 8029, 9016, 9995}
```

Next the fraction of the histogram occupied by “< 1000” is worked out. This is the selectivity. The histogram divides the range into equal frequency buckets, so all we have to do is locate the bucket that our value is in and count *part* of it and *all* of the ones before. The value 1000 is clearly in the second bucket (993–1997). Assuming a linear distribution of values inside each bucket, we can calculate the selectivity as:

```
selectivity = (1 + (1000 - bucket[2].min)/(bucket[2].max - bucket[2].min))/num_buckets
            = (1 + (1000 - 993)/(1997 - 993))/10
            = 0.100697
```

that is, one whole bucket plus a linear fraction of the second, divided by the number of buckets. The estimated number of rows can now be calculated as the product of the selectivity and the cardinality of `tenk1`:

```
rows = rel_cardinality * selectivity
      = 10000 * 0.100697
      = 1007   (rounding off)
```

Next let's consider an example with an equality condition in its `WHERE` clause:

```
EXPLAIN SELECT * FROM tenk1 WHERE stringu1 = 'CRAAAA';
```

QUERY PLAN

```
Seq Scan on tenk1   (cost=0.00..483.00 rows=30 width=244)
  Filter: (stringu1 = 'CRAAAA'::name)
```

Again the planner examines the `WHERE` clause condition and looks up the selectivity function for `=`, which is `eqsel`. For equality estimation the histogram is not useful; instead the list of *most common values* (MCVs) is used to determine the selectivity. Let's have a look at the MCVs, with some additional columns that will be useful later:

```
SELECT null_frac, n_distinct, most_common_vals, most_common_freqs FROM pg_stats
WHERE tablename='tenk1' AND attname='stringu1';
```

```
null_frac           | 0
n_distinct          | 676
most_common_vals    | {EJAAAA,BBAAAA,CRAAAA,FCAAAA,FEAAAA,GSAAAA,
JOAAAA,MCAAAA,NAAAAA,WGAAAA}
most_common_freqs   | {0.00333333,0.003,0.003,0.003,0.003,0.003,0.003,0.003,0.003,0.003}
```

Since `CRAAAA` appears in the list of MCVs, the selectivity is merely the corresponding entry in the list of most common frequencies (MCFs):

```
selectivity = mcf[3]
            = 0.003
```

As before, the estimated number of rows is just the product of this with the cardinality of `tenk1`:

```
rows = 10000 * 0.003
      = 30
```

Now consider the same query, but with a constant that is not in the MCV list:

```
EXPLAIN SELECT * FROM tenk1 WHERE stringu1 = 'xxx';
```

QUERY PLAN

```
Seq Scan on tenk1   (cost=0.00..483.00 rows=15 width=244)
  Filter: (stringu1 = 'xxx'::name)
```

This is quite a different problem: how to estimate the selectivity when the value is *not* in the MCV list. The approach is to use the fact that the value is not in the list, combined with the knowledge of the frequencies for all of the MCVs:

```
selectivity = (1 - sum(mcv_fregs))/(num_distinct - num_mcv)
            = (1 - (0.00333333 + 0.003 + 0.003 + 0.003 + 0.003 + 0.003 +
                    0.003 + 0.003 + 0.003 + 0.003))/(676 - 10)
            = 0.0014559
```

That is, add up all the frequencies for the MCVs and subtract them from one, then divide by the number of *other* distinct values. This amounts to assuming that the fraction of the column that is not any of the MCVs is evenly distributed among all the other distinct values. Notice that there are no null values so we don't have to worry about those (otherwise we'd subtract the null fraction from the numerator as well). The estimated number of rows is then calculated as usual:

```
rows = 10000 * 0.0014559
      = 15   (rounding off)
```

The previous example with `unique1 < 1000` was an oversimplification of what `scalarltsel` really does; now that we have seen an example of the use of MCVs, we can fill in some more detail. The example was correct as far as it went, because since `unique1` is a unique column it has no MCVs (obviously, no value is any more common than any other value). For a non-unique column, there will normally be both a histogram and an MCV list, and *the histogram does not include the portion of the column population represented by the MCVs*. We do things this way because it allows more precise estimation. In this situation `scalarltsel` directly applies the condition (e.g., “< 1000”) to each value of the MCV list, and adds up the frequencies of the MCVs for which the condition is true. This gives an exact estimate of the selectivity within the portion of the table that is MCVs. The histogram is then used in the same way as above to estimate the selectivity in the portion of the table that is not MCVs, and then the two numbers are combined to estimate the overall selectivity. For example, consider

```
EXPLAIN SELECT * FROM tenk1 WHERE stringu1 < 'IAAAAA';
```

QUERY PLAN

```
-----
Seq Scan on tenk1   (cost=0.00..483.00 rows=3077 width=244)
  Filter: (stringu1 < 'IAAAAA'::name)
```

We already saw the MCV information for `stringu1`, and here is its histogram:

```
SELECT histogram_bounds FROM pg_stats
WHERE tablename='tenk1' AND attname='stringu1';
```

histogram_bounds

```
-----
{AAAAAA,CQAAAA,FRAAAA,IBAAAA,KRAAAA,NFAAAA,PSAAAA,SGAAAA,VAAAAA,XLAAAA,ZZAAAA}
```

Checking the MCV list, we find that the condition `stringu1 < 'IAAAAA'` is satisfied by the first six entries and not the last four, so the selectivity within the MCV part of the population is

```
selectivity = sum(relevant mvfs)
            = 0.00333333 + 0.003 + 0.003 + 0.003 + 0.003 + 0.003
            = 0.01833333
```

Summing all the MCFs also tells us that the total fraction of the population represented by MCVs is 0.03033333, and therefore the fraction represented by the histogram is 0.96966667 (again, there are no nulls, else we'd have to exclude them here). We can see that the value `IAAAAA` falls nearly at the end of the third histogram bucket. Using some rather cheesy assumptions about the frequency of different characters, the planner arrives at the estimate 0.298387 for the portion of the histogram population that is less than `IAAAAA`. We then combine the estimates for the MCV and non-MCV populations:

```
selectivity = mcv_selectivity + histogram_selectivity * histogram_fraction
            = 0.01833333 + 0.298387 * 0.96966667
            = 0.307669

rows       = 10000 * 0.307669
            = 3077   (rounding off)
```

In this particular example, the correction from the MCV list is fairly small, because the column distribution is actually quite flat (the statistics showing these particular values as being more common than others are mostly due to sampling error). In a more typical case where some values are significantly more common than others, this complicated process gives a useful improvement in accuracy because the selectivity for the most common values is found exactly.

Now let's consider a case with more than one condition in the `WHERE` clause:

```
EXPLAIN SELECT * FROM tenk1 WHERE unique1 < 1000 AND stringu1 = 'xxx';
```

QUERY PLAN

```
-----
Bitmap Heap Scan on tenk1  (cost=23.80..396.91 rows=1 width=244)
  Recheck Cond: (unique1 < 1000)
  Filter: (stringu1 = 'xxx'::name)
    -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..23.80 rows=1007 width=0)
        Index Cond: (unique1 < 1000)
```

The planner assumes that the two conditions are independent, so that the individual selectivities of the clauses can be multiplied together:

```
selectivity = selectivity(unique1 < 1000) * selectivity(stringu1 = 'xxx')
              = 0.100697 * 0.0014559
              = 0.0001466

rows        = 10000 * 0.0001466
              = 1   (rounding off)
```

Notice that the number of rows estimated to be returned from the bitmap index scan reflects only the condition used with the index; this is important since it affects the cost estimate for the subsequent heap fetches.

Finally we will examine a query that involves a join:

```
EXPLAIN SELECT * FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 50 AND t1.unique2 = t2.unique2;
```

QUERY PLAN

```
-----
Nested Loop  (cost=4.64..456.23 rows=50 width=488)
  -> Bitmap Heap Scan on tenk1 t1  (cost=4.64..142.17 rows=50 width=244)
      Recheck Cond: (unique1 < 50)
      -> Bitmap Index Scan on tenk1_unique1  (cost=0.00..4.63 rows=50 width=0)
          Index Cond: (unique1 < 50)
  -> Index Scan using tenk2_unique2 on tenk2 t2  (cost=0.00..6.27 rows=1 width=244)
      Index Cond: (unique2 = t1.unique2)
```

The restriction on `tenk1, unique1 < 50`, is evaluated before the nested-loop join. This is handled analogously to the previous range example. This time the value 50 falls into the first bucket of the `unique1` histogram:

```
selectivity = (0 + (50 - bucket[1].min)/(bucket[1].max - bucket[1].min))/num_buckets
              = (0 + (50 - 0)/(993 - 0))/10
              = 0.005035

rows        = 10000 * 0.005035
              = 50   (rounding off)
```

The restriction for the join is `t2.unique2 = t1.unique2`. The operator is just our familiar `=`, however the selectivity function is obtained from the `oprjoin` column of `pg_operator`, and is `eqjoinsel`. `eqjoinsel` looks up the statistical information for both `tenk2` and `tenk1`:

```
SELECT tablename, null_frac, n_distinct, most_common_vals FROM pg_stats
```



```
WHERE tablename IN ('tenk1', 'tenk2') AND attname='unique2';
```

tablename	null_frac	n_distinct	most_common_vals
tenk1	0	-1	
tenk2	0	-1	

In this case there is no MCV information for `unique2` and all the values appear to be unique (`n_distinct = -1`), so we use an algorithm that relies on the row count estimates for both relations (`num_rows`, not shown, but "tenk") together with the column null fractions (zero for both):

```
selectivity = (1 - null_frac1) * (1 - null_frac2) / max(num_rows1, num_rows2)
             = (1 - 0) * (1 - 0) / max(10000, 10000)
             = 0.0001
```

This is, subtract the null fraction from one for each of the relations, and divide by the row count of the larger relation (this value does get scaled in the non-unique case). The number of rows that the join is likely to emit is calculated as the cardinality of the Cartesian product of the two inputs, multiplied by the selectivity:

```
rows = (outer_cardinality * inner_cardinality) * selectivity
      = (50 * 10000) * 0.0001
      = 50
```

Had there been MCV lists for the two columns, `eqjoinsel` would have used direct comparison of the MCV lists to determine the join selectivity within the part of the column populations represented by the MCVs. The estimate for the remainder of the populations follows the same approach shown here.

Notice that we showed `inner_cardinality` as 10000, that is, the unmodified size of `tenk2`. It might appear from inspection of the `EXPLAIN` output that the estimate of join rows comes from `50 * 1`, that is, the number of outer rows times the estimated number of rows obtained by each inner index scan on `tenk2`. But this is not the case: the join relation size is estimated before any particular join plan has been considered. If everything is working well then the two ways of estimating the join size will produce about the same answer, but due to round-off error and other factors they sometimes diverge significantly.

68.2. Multivariate Statistics Examples

68.2.1. Functional Dependencies

Multivariate correlation can be demonstrated with a very simple data set — a table with two columns, both containing the same values:

```
CREATE TABLE t (a INT, b INT);
INSERT INTO t SELECT i % 100, i % 100 FROM generate_series(1, 10000) s(i);
ANALYZE t;
```

As explained in [Section 14.2](#), the planner can determine cardinality of `t` using the number of pages and rows obtained from `pg_class`:

```
SELECT relpages, reltuples FROM pg_class WHERE relname = 't';
```

relpages	reltuples
45	10000

The data distribution is very simple; there are only 100 distinct values in each column, uniformly distributed.

The following example shows the result of estimating a `WHERE` condition on the `a` column:

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a = 1;
QUERY PLAN
```

```
Seq Scan on t (cost=0.00..170.00 rows=100 width=8) (actual rows=100 loops=1)
  Filter: (a = 1)
  Rows Removed by Filter: 9900
```

The planner examines the condition and determines the selectivity of this clause to be 1%. By comparing this estimate and the actual number of rows, we see that the estimate is very accurate (in fact exact, as the table is very small). Changing the `WHERE` condition to use the `b` column, an identical plan is generated. But observe what happens if we apply the same condition on both columns, combining them with `AND`:

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a = 1 AND b = 1;
                                QUERY PLAN
```

```
-----
Seq Scan on t (cost=0.00..195.00 rows=1 width=8) (actual rows=100 loops=1)
  Filter: ((a = 1) AND (b = 1))
  Rows Removed by Filter: 9900
```

The planner estimates the selectivity for each condition individually, arriving at the same 1% estimates as above. Then it assumes that the conditions are independent, and so it multiplies their selectivities, producing a final selectivity estimate of just 0.01%. This is a significant underestimate, as the actual number of rows matching the conditions (100) is two orders of magnitude higher.

This problem can be fixed by creating a statistics object that directs `ANALYZE` to calculate functional-dependency multivariate statistics on the two columns:

```
CREATE STATISTICS stts (dependencies) ON a, b FROM t;
ANALYZE t;
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a = 1 AND b = 1;
                                QUERY PLAN
```

```
-----
Seq Scan on t (cost=0.00..195.00 rows=100 width=8) (actual rows=100 loops=1)
  Filter: ((a = 1) AND (b = 1))
  Rows Removed by Filter: 9900
```

68.2.2. Multivariate N-Distinct Counts

A similar problem occurs with estimation of the cardinality of sets of multiple columns, such as the number of groups that would be generated by a `GROUP BY` clause. When `GROUP BY` lists a single column, the `n-distinct` estimate (which is visible as the estimated number of rows returned by the `HashAggregate` node) is very accurate:

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT COUNT(*) FROM t GROUP BY a;
                                QUERY PLAN
```

```
-----
HashAggregate (cost=195.00..196.00 rows=100 width=12) (actual rows=100 loops=1)
  Group Key: a
    -> Seq Scan on t (cost=0.00..145.00 rows=10000 width=4) (actual rows=10000
        loops=1)
```

But without multivariate statistics, the estimate for the number of groups in a query with two columns in `GROUP BY`, as in the following example, is off by an order of magnitude:

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT COUNT(*) FROM t GROUP BY a, b;
                                QUERY PLAN
```

```
-----
HashAggregate (cost=220.00..230.00 rows=1000 width=16) (actual rows=100 loops=1)
  Group Key: a, b
    -> Seq Scan on t (cost=0.00..145.00 rows=10000 width=8) (actual rows=10000
        loops=1)
```

By redefining the statistics object to include `n-distinct` counts for the two columns, the estimate is much improved:

```

DROP STATISTICS stts;
CREATE STATISTICS stts (dependencies, ndistinct) ON a, b FROM t;
ANALYZE t;
EXPLAIN (ANALYZE, TIMING OFF) SELECT COUNT(*) FROM t GROUP BY a, b;
                                QUERY PLAN
-----
HashAggregate  (cost=220.00..221.00 rows=100 width=16) (actual rows=100 loops=1)
  Group Key: a, b
    -> Seq Scan on t  (cost=0.00..145.00 rows=10000 width=8) (actual rows=10000
loops=1)

```

68.2.3. MCV Lists

As explained in [Section 68.2.1](#), functional dependencies are very cheap and efficient type of statistics, but their main limitation is their global nature (only tracking dependencies at the column level, not between individual column values).

This section introduces multivariate variant of MCV (most-common values) lists, a straightforward extension of the per-column statistics described in [Section 68.1](#). These statistics address the limitation by storing individual values, but it is naturally more expensive, both in terms of building the statistics in `ANALYZE`, storage and planning time.

Let's look at the query from [Section 68.2.1](#) again, but this time with a MCV list created on the same set of columns (be sure to drop the functional dependencies, to make sure the planner uses the newly created statistics).

```

DROP STATISTICS stts;
CREATE STATISTICS stts2 (mcv) ON a, b FROM t;
ANALYZE t;
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a = 1 AND b = 1;
                                QUERY PLAN
-----
Seq Scan on t  (cost=0.00..195.00 rows=100 width=8) (actual rows=100 loops=1)
  Filter: ((a = 1) AND (b = 1))
  Rows Removed by Filter: 9900

```

The estimate is as accurate as with the functional dependencies, mostly thanks to the table being fairly small and having a simple distribution with a low number of distinct values. Before looking at the second query, which was not handled by functional dependencies particularly well, let's inspect the MCV list a bit.

Inspecting the MCV list is possible using `pg_mcv_list_items` set-returning function.

```

SELECT m.* FROM pg_statistic_ext join pg_statistic_ext_data on (oid = stxoid),
       pg_mcv_list_items(stxmcv) m WHERE stxname = 'stts2';

```

index	values	nulls	frequency	base_frequency
0	{0, 0}	{f,f}	0.01	0.0001
1	{1, 1}	{f,f}	0.01	0.0001
...				
49	{49, 49}	{f,f}	0.01	0.0001
50	{50, 50}	{f,f}	0.01	0.0001
...				
97	{97, 97}	{f,f}	0.01	0.0001
98	{98, 98}	{f,f}	0.01	0.0001
99	{99, 99}	{f,f}	0.01	0.0001

(100 rows)

This confirms there are 100 distinct combinations in the two columns, and all of them are about equally likely (1% frequency for each one). The base frequency is the frequency computed from per-column

statistics, as if there were no multi-column statistics. Had there been any null values in either of the columns, this would be identified in the `nulls` column.

When estimating the selectivity, the planner applies all the conditions on items in the MCV list, and then sums the frequencies of the matching ones.

Compared to functional dependencies, MCV lists have two major advantages. Firstly, the list stores actual values, making it possible to decide which combinations are compatible.

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a = 1 AND b = 10;
                                QUERY PLAN
```

```
-----
Seq Scan on t   (cost=0.00..195.00 rows=1 width=8) (actual rows=0 loops=1)
  Filter: ((a = 1) AND (b = 10))
  Rows Removed by Filter: 10000
```

Secondly, MCV lists handle a wider range of clause types, not just equality clauses like functional dependencies. For example, consider the following range query for the same table:

```
EXPLAIN (ANALYZE, TIMING OFF) SELECT * FROM t WHERE a <= 49 AND b > 49;
                                QUERY PLAN
```

```
-----
Seq Scan on t   (cost=0.00..195.00 rows=1 width=8) (actual rows=0 loops=1)
  Filter: ((a <= 49) AND (b > 49))
  Rows Removed by Filter: 10000
```

68.3. Planner Statistics and Security

Access to the table `pg_statistic` is restricted to superusers, so that ordinary users cannot learn about the contents of the tables of other users from it. Some selectivity estimation functions will use a user-provided operator (either the operator appearing in the query or a related operator) to analyze the stored statistics. For example, in order to determine whether a stored most common value is applicable, the selectivity estimator will have to run the appropriate `=` operator to compare the constant in the query to the stored value. Thus the data in `pg_statistic` is potentially passed to user-defined operators. An appropriately crafted operator can intentionally leak the passed operands (for example, by logging them or writing them to a different table), or accidentally leak them by showing their values in error messages, in either case possibly exposing data from `pg_statistic` to a user who should not be able to see it.

In order to prevent this, the following applies to all built-in selectivity estimation functions. When planning a query, in order to be able to use stored statistics, the current user must either have `SELECT` privilege on the table or the involved columns, or the operator used must be `LEAKPROOF` (more accurately, the function that the operator is based on). If not, then the selectivity estimator will behave as if no statistics are available, and the planner will proceed with default or fall-back assumptions.

If a user does not have the required privilege on the table or columns, then in many cases the query will ultimately receive a permission-denied error, in which case this mechanism is invisible in practice. But if the user is reading from a security-barrier view, then the planner might wish to check the statistics of an underlying table that is otherwise inaccessible to the user. In that case, the operator should be leak-proof or the statistics will not be used. There is no direct feedback about that, except that the plan might be suboptimal. If one suspects that this is the case, one could try running the query as a more privileged user, to see if a different plan results.

This restriction applies only to cases where the planner would need to execute a user-defined operator on one or more values from `pg_statistic`. Thus the planner is permitted to use generic statistical information, such as the fraction of null values or the number of distinct values in a column, regardless of access privileges.

Selectivity estimation functions contained in third-party extensions that potentially operate on statistics with user-defined operators should follow the same security rules.

Chapter 69. Backup Manifest Format

The backup manifest is primarily intended to permit the backup to be verified using [pg_verifybackup](#). However, it is also possible for other tools to read the backup manifest file and use the information contained therein for their own purposes. To that end, this chapter describes the format of the backup manifest file.

A backup manifest is a JSON document encoded as UTF-8. (Although in general JSON documents are required to be Unicode, Postgres Pro permits the `json` and `jsonb` data types to be used with any supported server encoding. There is no similar exception for backup manifests.) The JSON document is always an object; the keys that are present in this object are described in the next section.

69.1. Backup Manifest Top-level Object

The backup manifest JSON document contains the following keys.

`PostgreSQL-Backup-Manifest-Version`

The associated value is an integer. Beginning in Postgres Pro 17, it is 2; in older versions, it is 1.

`System-Identifier`

The database system identifier of the Postgres Pro instance where the backup was taken. This field is present only when `PostgreSQL-Backup-Manifest-Version` is 2.

`Files`

The associated value is always a list of objects, each describing one file that is present in the backup. No entries are present in this list for the WAL files that are needed in order to use the backup, or for the backup manifest itself. The structure of each object in the list is described in [Section 69.2](#).

`WAL-Ranges`

The associated value is always a list of objects, each describing a range of WAL records that must be readable from a particular timeline in order to make use of the backup. The structure of these objects is further described in [Section 69.3](#).

`Manifest-Checksum`

This key is always present on the last line of the backup manifest file. The associated value is a SHA256 checksum of all the preceding lines. We use a fixed checksum method here to make it possible for clients to do incremental parsing of the manifest. While a SHA256 checksum is significantly more expensive than a CRC32C checksum, the manifest should normally be small enough that the extra computation won't matter very much.

69.2. Backup Manifest File Object

The object which describes a single file contains either a `Path` key or an `Encoded-Path` key. Normally, the `Path` key will be present. The associated string value is the path of the file relative to the root of the backup directory. Files located in a user-defined tablespace will have paths whose first two components are `pg_tblspc` and the OID of the tablespace. If the path is not a string that is legal in UTF-8, or if the user requests that encoded paths be used for all files, then the `Encoded-Path` key will be present instead. This stores the same data, but it is encoded as a string of hexadecimal digits. Each pair of hexadecimal digits in the string represents a single octet.

The following two keys are always present:

`Size`

The expected size of this file, as an integer.

Last-Modified

The last modification time of the file as reported by the server at the time of the backup. Unlike the other fields stored in the backup, this field is not used by [pg_verifybackup](#). It is included only for informational purposes.

If the backup was taken with file checksums enabled, the following keys will be present:

Checksum-Algorithm

The checksum algorithm used to compute a checksum for this file. Currently, this will be the same for every file in the backup manifest, but this may change in future releases. At present, the supported checksum algorithms are CRC32C, SHA224, SHA256, SHA384, and SHA512.

Checksum

The checksum computed for this file, stored as a series of hexadecimal characters, two for each byte of the checksum.

69.3. Backup Manifest WAL Range Object

The object which describes a WAL range always has three keys:

Timeline

The timeline for this range of WAL records, as an integer.

Start-LSN

The LSN at which replay must begin on the indicated timeline in order to make use of this backup. The LSN is stored in the format normally used by Postgres Pro; that is, it is a string consisting of two strings of hexadecimal characters, each with a length of between 1 and 8, separated by a slash.

End-LSN

The earliest LSN at which replay on the indicated timeline may end when making use of this backup. This is stored in the same format as `Start-LSN`.

Ordinarily, there will be only a single WAL range. However, if a backup is taken from a standby which switches timelines during the backup due to an upstream promotion, it is possible for multiple ranges to be present, each with a different timeline. There will never be multiple WAL ranges present for the same timeline.

Chapter 70. Distributed System Internals

The Postgres Pro Shardman software comprises these main components: PostgreSQL core with additional features, shardman extension, management services and utilities. This section considers Postgres Pro Shardman cluster as a group of PostgreSQL instances or shards. Each shard may also have several standby nodes and to emphasize this the term replication group is used. The support for highly available configurations is currently done on the level of tools and services and will be covered in the Management section.

70.1. Table Types

In a distributed database managed by Postgres Pro Shardman the following special table types are used: sharded tables and global tables.

70.1.1. Sharded Tables

Sharded tables are just usual PostgreSQL partitioned tables where a few partitions, making up a *shard*, are regular local tables and the other partitions are foreign tables available from remote servers via [postgres_fdw](#). Sharded tables are registered in the `shardman.sharded_tables` dictionary. Use the `CREATE TABLE` statement with the `distributed_by` parameter to create a sharded table. Several sharded tables can be created as colocated. This means that they have the same number of partitions and that their partitions corresponding to the same sharding key should reside together. During a rebalance, Postgres Pro Shardman management utilities ensure that corresponding partitions of colocated tables are moved to the same node. (Such a rebalance happens, for example, when a new node is added to the cluster). Colocation is necessary to ensure that joins of several tables are propagated to the node where the actual data resides. To define one sharded table colocated with another one, first, create one table and then use the `colocate_with` parameter of the `CREATE TABLE` statement while creating the second table. Chains of colocated tables are not supported, all related tables should be marked as colocated to one of the tables instead. Note that `colocate_with` property is symmetric and transitive.

70.1.1.1. Partitions

A sharded table consists of several partitions. Some of them are regular tables, and others are foreign tables. By default, the number of partitions is determined by the `shardman.num_parts` parameter, but it can be overwritten by the `num_parts` `CREATE TABLE` parameter. Most of DDL operations are restricted on partitions of a sharded table. You should modify the parent table instead.

The number of partitions in a sharded table is defined when it is created and cannot be changed afterwards. When new nodes are added to the cluster, some partitions are moved from existing nodes to the new ones to balance the load. So, to allow scaling of clusters, the initial number of partitions should be high enough, but not too high since an extremely large number of partitions significantly slows down query planning. For example, if you expect the number of nodes in your cluster to grow by 4 times at a maximum, create sharded tables with the number of partitions equal to $4 * N$, where N is the number of nodes. A cluster becomes unable to scale when the number of cluster nodes reaches the number of partitions in the sharded table with the minimal number of them.

70.1.1.2. Subpartitions

Partitions of a sharded table can be partitioned by range. In this case, each partition of a sharded table is a partitioned table consisting only of regular or only of foreign subpartitions. All subpartitions of a partition are located on the same node. Use the `partition_by` `CREATE TABLE` parameter to specify a column that should be used as a subpartition key column and the `partition_bounds` parameter to set bounds of the second-level table partitions. New subpartitions can be added or removed from a table as necessary. So you can omit the `partition_bounds` parameter during table creation and create partitions later using the `shardman.create_subpart()` function. Other subpartition management functions allow you to drop, detach or attach subpartitions of a sharded table. Subpartition management is cluster-wide.

70.1.2. Global Tables

Global tables are available to all nodes of a cluster. Now a global table is a set of regular tables synchronized by triggers. The main use case for a global table is to store a relatively rarely updated set of data that is used by all cluster nodes. When a sharded table is joined to a global table, joins between sharded table partitions and the global table can be performed on nodes where individual partitions reside. The implementation of trigger-based replication requires a non-deferrable primary key on a global table to be defined. Currently when a global table is modified, an after-statement trigger fires and propagates changes to other nodes of the cluster via foreign tables. When new nodes are added to a cluster, global table data is transferred to the new nodes via logical replication. When some nodes are removed from a cluster, global tables get locked for writes for a brief time. Use the `global CREATE TABLE` parameter to create a global table. Global tables are registered in the `shardman.global_tables` dictionary. Partitioned global tables are not supported.

70.1.3. Distributed DDL

Postgres Pro Shardman extension allows creating several kinds of global objects. These are sharded and global tables, roles, sequences and tablespaces. The list of operations allowed on global objects is limited particularly to protect consistency of a global schema. For the same reason, most operations on global objects are cluster-wide. The list of cluster-wide operations includes:

- `CREATE` for sharded and global tables, global roles and tablespaces or indexes on sharded or global tables.
- `DROP` for sharded and global tables, global roles and tablespaces or indexes on sharded or global tables.
- `ALTER TABLE` for sharded and global tables.
- `ALTER TABLESPACE` for global tablespaces.
- `ALTER ROLE` for global roles.
- `RENAME` for sharded and global tables or indexes on them.
- `SET CONSTRAINTS ALL` inside a transaction block.

These configuration settings control execution of the distributed DDL: [shardman.broadcast_ddl](#) and [shardman.sync_schema](#). The first one can be used for a cluster-wide broadcast of all regular DDL operations (for example, creating schemas or functions). The second one controls broadcasting of statements related to global objects and should never be turned off without consulting the Postgres Pro Shardman support team.

70.2. Silk

70.2.1. Concept

Silk (Postgres Pro Shardman InterLinK) is a transport feature used for inter-cluster communications. It is injected at the point where `postgres_fdw` decides to transmit deparsed piece of query through `libpq` connection to the remote node, replacing `libpq` connection with itself. It is designed to decrease the count of idle `postgres_fdw` connections during transaction execution, minimize latency and boost overall throughput.

Silk implementation uses several background processes. The main routing/multiplexing process (one per PostgreSQL instance), called `silkroad`, and a bunch of background workers, called `silkworms`. While `postgres_fdw` uses `libpq`, it spawns multiple `libpq` connections from each backend to the remote node (where multiple backend processes are spawned accordingly). But if `silk` replaces `libpq` - every `silkroad` process is connected to only one remote `silkroad`. In this scheme, remote `silkworms` play the role of remote backends otherwise spawned by `postgres_fdw`.

`Silkroad` wires local backend with remote node's workers this way:

1. Backend process uses regular `postgres_fdw` API to access remote data as usual. But `postgres_fdw`, when silk is enabled, writes the query into shared memory queue instead of `libpq` connection;
2. `Silkroad` process parses incoming shared memory queue from that backend and routes the message to appropriate network connection with remote `silkroad` process.
3. Remote `silkroad` process grabs incoming message from network and (if it is a new one) redirects it to available worker's shared memory queue (or in a special "unassigned jobs" queue if all of the workers are busy).
4. At last, remote worker gets the message through its shared memory queue, executes it and sends back the result tuples (or an error) the same way.

`Silkroad` acts here like a common network switch, tossing packets between backend's shared memory and appropriate network socket. It knows nothing about content of a message relying only on the message header.

70.2.2. Event Loop

`Silkroad` process runs an event loop powered by the `libev` library. Each backend's shared memory queue is exposed at the event loop with the `eventfd` descriptor, and each network connection - with a socket descriptor.

During startup, the backend registers itself (its `eventfd` descriptors) at a local `silkroad` process. `Silkroad` responds by specifying which memory segments to use for the backend's message queue. From this moment `silkroad` will respond to events from the queue associated with this backend. Network connections between local and remote `silkroads` will be established at once on the first request from the backend to the remote node and stay alive until both of participants (`silkroad` processes) exist.

70.2.3. Routing and Multiplexing

For each subquery, we expect a subset of tuples, and therefore represent the interaction within the subquery as a bidirectional data stream. `Silkroad` uses an internal routing table to register these streams. A unique stream ID (within the Postgres Pro Shardman cluster) is formed as a pair of "origin node address, target node address" and a locally (within the node) unique number. Each particular subquery from a backend to remote nodes will be registered by `silkroad` as such a stream. So, any backend can be associated with many streams at the time.

When a local `silkroad` process got a message with a new stream ID from backend, it registers it in a local routing table and then redirects this message to an appropriate socket. If the connection with the remote `silkroad` does not exist, it is established using a handshake procedure. The original message that initiated a handshake is placed into a special internal buffer until the handshake succeeds. The remote `silkroad` process receiving a packet with the new ID registers it in its own table, then assigns a `silkworm` worker from a pool of available workers and places the message into the worker's shared memory queue. If all of the `silkworm` workers are busy at the moment, the message will be postponed, i.e., placed into a special "unassigned jobs queue" (note that the configuration parameter `shardman.silk_unassigned_job_queue_size` is 1024). If there is no free space in the queue, an error message will be generated and sent back to the source backend. A job from this queue will be assigned later to the first available worker when it gets rid of the previous job.

When the worker got a new "job", it executes it through `SPI` subsystem, organizing result tuples into batches and sends them back through shared memory to the local `silkroad` process. The rest is trivial due to the whole route is known. The last resulting packet with tuples in a stream is marked as "closing". It is an order to `silkroads` to wipe out this route from their tables.

Note that backend and remote workers stay "subscribed" to their streams until they are explicitly closed. So the backend has the opportunity to send an abort message or notify the remote worker to prematurely close the transaction. And it makes it possible to discard obsolete data packets, possibly from previous aborted transactions.

To observe the current state of the `silkroad` multiplexer process, Silk diagnostics views are available.

70.2.4. Error Handling and Route Integrity

Besides the routing table `silkroad` tracks endpoints (backends and network connections) that were involved in some particular stream. So when some connection is closed, all the involved backends (and/or workers) will be notified of that event with a special error message, and all routes/streams related to this connection will be dismissed. The same way, if the backend crashes, its shared memory queue become detached and `silkroad` reacts by sending error messages to remote participants of every stream related to the crashed backend. So remote workers are not left doing useless work when the requester has already died.

70.2.5. Data Transmitting/batching/splitting Oversized Tuples

The resulting tuples are transmitted by `silkworm` in a native binary mode. Tuples with `external` storage attribute will be deTOASTed, but those that were compressed stay compressed.

Small tuples will be organized in batches (about 256k). Big tuples will be cut into pieces by the sender and assembled into a whole by the receiving backend.

70.2.6. Streams Flow Control

It may happen that when the next message is received from a backend, it will not fit the target network buffer. Or the message received from the network does not fit into the target shared memory queue. In such a case, the stream that caused this situation will be “suspended”. This means that the `silkroad` pauses the reaction to events from the source endpoint (connection or backend) until the target endpoint drains their messages. The rest backends and connections not affected by this route are kept working. Receiving modules of backends are designed to minimize these situations. The backend periodically checks and drains the incoming queue even when the plan executor is busy processing other plan nodes. Received tuples are stored in backend's tuplestores according the plan nodes until the executor requests the next tuple for a particular plan node execution.

When enough space is freed on the target queue, the suspended stream gets resumed, endpoint's events get unblocked and the process of receiving and sorting packets continues.

70.2.7. Implementation details

70.2.7.1. State Transferring and CSNs

When `postgres_fdw` works over Silk transport, only one connection between `silkroad` routing daemons is used to transfer user requests to `silkworm` workers and get their responses. Each request contains a transaction state, a replication group ID of the node where the request is formed (coordinator), a query itself and query parameters (if present). A response is either an error response message with a specific error message and error code or a bunch of tuples followed by “end of tuples” message. This means that `silkworm` has to switch to the transaction state coming with the request prior to executing the request.

For now, Silk transport is used only for read-only `SELECT` queries. All modifying requests are processed via a usual `libpq` connection and handled mostly as all other DML requests in PostgreSQL `postgres_fdw`. The only distinction is that when a DML request is processed by `postgres_fdw`, it resets the saved transaction state for the connection cache entry corresponding to the connection where this request is sent. Also a read-only flag is set to false for such a connection cache entry. When a request is sent over Silk transport, Postgres Pro Shardman extension asks for the transaction state for a pair of `serverid` and `userid` from `postgres_fdw`. If such a connection cache entry is found in the `postgres_fdw` connection cache, it is not a read-only cache entry and transaction state is present in this entry, the state is returned. If it is not present, `postgres_fdw` retrieves a full transaction state from the remote server, saves it in the connection cache entry and returns to the Postgres Pro Shardman extension.

The full transaction state is similar to the parallel worker transaction state and contains:

- information related to the current user (`uid`, `username`)
- `pid` of the current backend

- transaction start timestamp
- current snapshot CSN
- flags indicating that invalidation messages are present
- set of configuration settings, set in current session (including the application name)
- backend private state:
 - array of ComboCIDs
 - internal transaction state (full transaction ID, isolation level, current command ID, etc.)
 - information about reindexed indexes

If the connection is not found in the `postgres_fdw` connection cache (i.e., it is a new connection) or the entry in the connection cache is marked as read-only, only these characteristics form the transaction state:

- information related to the current user (username)
- transaction start timestamp
- current snapshot CSN
- flags indicating that invalidation messages are present
- set of configuration settings, set in current session (including application name)

Using such transaction states, `silkworm` can attach to a running transaction or start a new read-only transaction with the provided snapshot CSN and retrieve the result.

Note that the full transaction state can be imported only on the server that exported it. Also note that due to this transaction state transferring method, you cannot use Silk transport without enabling CSN snapshots.

70.2.7.2. Integration with Asynchronous FDW Engine

In the [Section 53.2.2](#), asynchronous `ForeignScan` plan nodes were presented as a way to optimize data retrieval from multiple hosts while these plan nodes were located under a single `Append` node. In the standard PostgreSQL architecture, the execution of `ForeignScan` plan nodes is implemented using the network protocol based on `libpq`. To improve the system performance during data transfer and reduce resource consumption, Postgres Pro Shardman employs a different method for exchanging data with remote hosts. The mechanism for executing `ForeignScan` nodes is implemented using the Silk protocol.

To incorporate Silk transport into the asynchronous executor, modifications were made to the `postgres_fdw` extension. A pluggable transport was implemented as a set of interface functions included as part of the Postgres Pro Shardman extension. During execution of callbacks that interact with remote hosts, these functions are called by the `postgres_fdw` extension. The pluggable Silk transport is activated if the Postgres Pro Shardman extension is preloaded and if the foreign server has the attribute `extended_features` (applicable for any FDW server in the Postgres Pro Shardman cluster). For all other cases, the `postgres_fdw` extension uses the standard exchange protocol based on `libpq`.

To disable the pluggable Silk transport for query execution, it is necessary to set the `query_engine_mode` configuration parameter to the value of `ENGINE_NONE`.

In the current implementation, the pluggable Silk transport is only used for read-only queries, specifically during the execution of the `ForeignScan` node. The standard exchange protocol based on `libpq` is used for modifying queries.

When receiving query execution result rows using the Silk transport, the data is stored in a `TupleStoreState` storage as a complete result set, which is the same size as that returned by the remote host. The `TupleStoreState` is implemented as a data structure that can spill data to the disk in case of memory shortage. If the remote host returns a large result set, it does not lead to an out-of-memory

(OOM) condition. Once the result set is received in the `TupleStoreState`, the data is copied into the `ForeignScan` executor's in-memory buffer. The size of this buffer is defined by the `fetch_size` attribute of the foreign server. The default value of 50000 rows can be adjusted to find a balance between the performance (number of `ForeignScan` node calls) and memory consumption.

Utilizing the pluggable Silk transport for the asynchronous FDW engine results in an increase of the network exchange performance and a reduction of the system resource consumption due to better utilization of system resources, including the number of network connections.

70.2.8. Multiplexor Diagnostics Views

Views in this section provide various information related to Silk multiplexing. See [Section 70.2.3](#) for details of `silkroad` multiplexing process.

70.2.8.1. `shardman.silk_routes`

The `shardman.silk_routes` view displays the current snapshot of the multiplexer routing table. The columns of the view are shown in [Table 70.1](#).

Table 70.1. `shardman.silk_routes` Columns

Name	Type	Description
<code>hashvalue</code>	<code>integer</code>	Internal unique route identifier. Can be used to join with other Silk diagnostics views.
<code>origin_ip</code>	<code>inet</code>	IP address of the source node, which generated this route
<code>origin_port</code>	<code>int2</code>	External TCP connection port of the source node, which generated this route
<code>channel_id</code>	<code>integer</code>	Route sequential number within the node that generated this route. <code>channel_id</code> is unique for the pair <code>origin_ip</code> + <code>origin_port</code> . This pair is a unique node identifier within the Postgres Pro Shardman cluster and hence the <code>origin_ip</code> + <code>origin_port</code> + <code>channel_id</code> tuple is a unique route identifier within the Postgres Pro Shardman cluster.
<code>from_cn</code>	<code>integer</code>	Connect index in the <code>shardman.silk_connects</code> view for incoming routes, that is, not generated by this node, and -1 for routes generated by this node.
<code>backend_id</code>	<code>integer</code>	ID of the local process that is currently using this route: either the ID of the backend that generated this route or the ID of the <code>silkworm</code> worker assigned to this route. Equals -1 for queued incoming routes that have not been assigned a worker yet.
<code>pending_queue_bytes</code>	<code>bigint</code>	Size of the queue of delayed messages (awaiting a free work-

Name	Type	Description
		er) for this route, in bytes. This value is only meaningful for incoming routes of each node that are not assigned to a worker yet.
pending_queue_messages	bigint	Number of messages in the queue of delayed messages (awaiting a free worker) for this route. This value is only meaningful for incoming routes of each node that are not assigned to a worker yet.
connects	integer[]	List of indexes of connects that are currently using this route.

70.2.8.2. shardman.silk_connects

The `shardman.silk_connects` view displays the current list of multiplexer connects. The columns of the view are shown in [Table 70.2](#).

Table 70.2. shardman.silk_connects Columns

Name	Type	Description
cn_index	integer	Unique connect index
reg_ip	inet	“Registration” IP address of the node with which the connection is established. See Notes for details.
reg_port	int2	“Registration” TCP port of the node with which the connection is established. See Notes for details.
read_ev_active	boolean	<code>true</code> if the multiplexer is ready to receive data to the incoming queue. See Notes for details.
write_ev_active	boolean	<code>true</code> if the multiplexer filled the queue of non-sent messages and is waiting for it to get free. See Notes for details.
is_outgoing	boolean	<code>true</code> if the connection is outgoing, that is, created by <code>connect</code> , and <code>false</code> for incoming connects, that is, created by <code>accept</code> . Only used during the handshaking.
state	text	Current state of the connect: <code>connected</code> — if the connection is established, <code>in progress</code> — if the client has already connected, but handshaking has not happened yet, <code>free</code> — if the client has already disconnected, but the connect structure for the disconnected client has not been destroyed yet.

Name	Type	Description
pending_queue_bytes	bigint	Size of the queue of non-sent messages for this connect, in bytes
pending_queue_messages	bigint	Number of messages in the queue of non-sent messages for this connect
blocked_by_backend	integer	ID of the backend that blocked this connect
blocks_backends	integer[]	List of IDs of backends that are blocked by this connect
routes	integer[]	List of unique IDs of routes that use this connect
elapsed_time_write	bigint	Time from the last writing event of a connect
elapsed_time_read	bigint	Time from the last reading event of a connect

70.2.8.3. shardman.silk_backends

The `shardman.silk_backends` view displays the current list of processes of two kinds: backends that serve client connections and `silkworm` multiplexer workers, which interact with the multiplexer. The columns of the view are shown in [Table 70.3](#).

Table 70.3. shardman.silk_backends Columns

Name	Type	Description
backend_id	integer	Unique backend/worker identifier
pid	integer	OS process ID
attached	boolean	Value is true if backend is attached to multiplexer, false otherwise
read_ev_active	boolean	true if the backend/worker is ready to receive data to the incoming queue. See Notes for details.
write_ev_active	boolean	true if the backend/worker filled the queue of non-sent messages and is waiting for it to get free. See Notes for details.
is_worker	boolean	true if this process is a <code>silkworm</code> multiplexer worker and false otherwise
pending_queue_bytes	bigint	Size of the queue of messages being sent to this backend/worker, in bytes
pending_queue_messages	bigint	Number of messages in the queue of messages being sent to this backend/worker
blocked_by_connect	integer	Index of the connect that blocks this backend/worker

Name	Type	Description
blocks_connects	integer[]	List of indexes of connects that are blocked by this backend/worker
routes	integer[]	List of unique IDs of routes that are used by this backend/worker
in_queue_used	bigint	Number of queued data bytes in the incoming queue in the shared memory between the backend and multiplexer
out_queue_used	bigint	Number of queued data bytes in the outgoing queue in the shared memory between the backend and multiplexer
elapsed_time_write	bigint	Time from the last writing event of a backend
elapsed_time_read	bigint	Time from the last reading event of backend

70.2.8.4. shardman.silk_routing

The `shardman.silk_routing` view displays the results of the `shardman.silk_routing` function. [Table 70.4](#).

Table 70.4. shardman.silk_routing Columns

Name	Type	Description
hashvalue	integer	Internal unique route identifier
origin_ip	inet	IP address of the node that generated this route
origin_port	int2	External TCP connection port of the source node that generated this route
channel_id	integer	Route sequential number within the node that generated this route
is_outgoing	boolean	<code>true</code> if this route was produced by the outgoing network connection, <code>false</code> if it was produced by the incoming network connection.
pending_queue_bytes	bigint	Pending queue size, in bytes
pending_queue_messages	bigint	Number of pending queue messages
backend_id	integer	ID of the local process that is currently using this route: either the ID of the backend that generated this route or the ID of the <code>silkworm</code> worker assigned to this route. Equals -1 for queued incoming routes that have not been assigned a worker yet.

Name	Type	Description
backend_pid	integer	Returns the process ID of the server process attached to the current session
attached	boolean	Value is true if backend is attached to multiplexer, false otherwise
backend_rd_active	boolean	true if the backend/worker is ready to receive data to the incoming queue. See Notes for details.
backend_wr_active	boolean	true if the backend/worker filled the queue of non-sent messages and is waiting for it to get free. See Notes for details.
is_worker	boolean	true if this process is a silkworm multiplexer worker and false otherwise
backend_blocked_by_cn	integer	Index of the connect that blocks this backend/worker
blocks_connects	integer[]	List of indexes of connects that are blocked by this backend/worker
in_queue_used	bigint	Number of queued data bytes in the incoming queue in the shared memory between the backend and multiplexer
out_queue_used	bigint	Number of queued data bytes in the outgoing queue in the shared memory between the backend and multiplexer
connect_id	integer	Unique connect index
reg_ip	inet	“Registration” IP address of the node with which the connection is established
reg_port	int2	“Registration” TCP port of the node with which the connection is established
connect_rd_active	boolean	true if the multiplexer is ready to receive data to the incoming queue
connect_wr_active	boolean	true if the multiplexer filled the queue of non-sent messages and is waiting for it to get free
connect_is_outgoing	boolean	true if the connection is outgoing, that is, created by connect, and false for incoming connects, that is, created by accept. Only used during the handshaking.

Name	Type	Description
connect_state	text	Current state of the connect: <code>connected</code> — if the connection is established, <code>in progress</code> — if the client has already connected, but handshaking has not happened yet, <code>free</code> — if the client has already disconnected, but the connect structure for the disconnected client has not been destroyed yet
connect_outgoing_queue_bytes	bigint	Size of the queue of non-sent messages for this connect, in bytes
connect_outgoing_queue_messages	bigint	Number of messages in the queue of non-sent messages for this connect
connect_blocked_by_bk	integer	ID of the backend that blocked this connect
blocks_backends	integer[]	List of IDs of backends that are blocked by this connect
connect_elapsed_time_write	bigint	Time from the last writing event of a connect
connect_elapsed_time_read	bigint	Time from the last reading event of a connect
backend_elapsed_time_write	bigint	Time from the last writing event of a backend
backend_elapsed_time_read	bigint	Time from the last reading event of a backend

70.2.8.5. shardman.silk_pending_jobs

The `shardman.silk_pending_jobs` view displays the current list of routes in the queue of delayed multiplexer jobs, that is, jobs that are not assigned to workers yet. The columns of the view are shown in [Table 70.5](#).

Table 70.5. shardman.silk_pending_jobs Columns

Name	Type	Description
hashvalue	integer	Internal unique route identifier
origin_ip	inet	IP address of the node that generated this route
origin_port	int2	TCP connection port of the node that generated this route
channel_id	integer	Route sequential number within the node that generated this route
query	text	The first queued message
pending_queue_bytes	bigint	Pending queue size, in bytes
pending_queue_messages	bigint	Number of pending queue messages

70.2.8.6. shardman.silk_statinfo

The `shardman.silk_statinfo` view displays the current multiplexer state information. The columns of the view are shown in [Table 70.6](#).

Table 70.6. `shardman.silk_statinfo` Columns

Name	Type	Description
<code>pid</code>	<code>integer</code>	<code>silkroad</code> process ID
<code>started_at</code>	timestamp with time zone	Time when the <code>silkroad</code> back-end was started.
<code>transferred_bytes</code>	<code>json</code>	JSON object of key value pairs, where the key is the name of the message type, and the value is total number of bytes sent for the message types with at least one message sent
<code>transferred_pkts</code>	<code>json</code>	JSON object of key value pairs, where the key is the name of the message type, and the value is the total number of sent messages for the message types with at least one message sent
<code>transferred_max</code>	<code>json</code>	JSON object of key value pairs, where the key is the name of the message type, and the value is the maximum size of a message for the message types with at least one message sent
<code>memcxt_dpg_allocated</code>	<code>bigint</code>	The <code>mem_allocated</code> value of the process in <code>DPGMemoryContext</code>
<code>memcxt_top_allocated</code>	<code>bigint</code>	The <code>mem_allocated</code> value of the process in <code>TopMemoryContext</code>
<code>read_efd_max</code>	<code>bigint</code>	Maximum reading time of the <code>eventfd</code> since reset
<code>write_efd_max</code>	<code>bigint</code>	Maximum writing time of the <code>eventfd</code> since reset
<code>read_efd_total</code>	<code>bigint</code>	Total reading time of the <code>eventfd</code> since reset
<code>write_efd_total</code>	<code>bigint</code>	Total writing time of the <code>eventfd</code> since reset
<code>read_efd_count</code>	<code>bigint</code>	Total number of reading events of the <code>eventfd</code> since reset
<code>write_efd_count</code>	<code>bigint</code>	Total number of writing events of the <code>eventfd</code> since reset
<code>sort_time_max</code>	<code>bigint</code>	Maximum time of sorting operations with the <code>silk_flow_control</code> enabled (any value other than <code>none</code>)
<code>sort_time_total</code>	<code>bigint</code>	Total time of sorting operations with the <code>silk_flow_control</code> enabled (any value other than <code>none</code>)

Name	Type	Description
sort_time_count	bigint	Total number of the sorting operations with the <code>silk_flow_control</code> enabled (any value other than none)

Note that `read_efd_max`, `write_efd_max`, `read_efd_total`, `write_efd_total`, `read_efd_count`, `write_efd_count`, `sort_time_max`, `sort_time_total`, and `sort_time_count` are only calculated if the `shardman.silk_track_time` configuration parameter is enabled.

70.2.8.7. `shardman.silk_state`

The `shardman.silk_state` view displays the current `silkroad` process state. The columns of the view are shown in Table 70.7.

Table 70.7. `shardman.silk_state` Columns

Name	Type	Description
state	text	State of the <code>silkroad</code> process

70.2.8.8. Notes

`reg_ip` and `reg_port` values are not actual network addresses, but the addresses by which the multiplexer accesses the node. They are determined during a handshake between multiplexer nodes and are equal to the corresponding parameters of an appropriate server in the `pg_foreign_server` table.

All the `read_ev_active` values are `true` and all the `write_ev_active` values are `false` when the multiplexer is in the `idle` state.

70.2.8.9. Global Views

Postgres Pro Shardman has a list of global views based on the corresponding local views. The definition of global view columns is the same as in its corresponding local view. Fetching from a global view returns a union of rows from the corresponding local views. The rows are fetched from each of their cluster nodes. Another difference is that the global views have an added column `rgid`. The `rgid` value shows the replication group ID of the cluster node from which a row is fetched.

Below is the list of the global views with links to their corresponding local views:

Table 70.8. Silk-related global and local views

Global view	Local view	Description
<code>shardman.gv_silk_routes</code>	shardman.silk_routes	One row showing the current snapshot of the multiplexer routing table.
<code>shardman.gv_silk_connects</code>	shardman.silk_connects	One row showing the current list of multiplexer connects.
<code>shardman.gv_silk_backends</code>	shardman.silk_backends	One row showing the current list of processes of two kinds: backends that serve client connections and <code>silkworm</code> multiplexer workers, which interact with the multiplexer.
<code>shardman.gv_silk_pending_jobs</code>	shardman.silk_pending_jobs	One row showing the current list of routes in the queue of multiplexer jobs that are not assigned to workers yet.

Global view	Local view	Description
<code>shardman.gv_silk_routing</code>	<code>shardman.silk_routing</code>	One row showing the results of the <code>shardman.silk_routing</code> function.

70.2.9. Functions

```
shardman.silk_statinfo_reset()
```

Resets the values of the metrics with prefix `transferred_` and time-based metrics (with prefixes `read_efd_`, `write_efd_`, and `sort_time_`) in the `shardman.silk_statinfo` view.

```
shardman.silk_routing()
```

Retrieves the results of the multiplexer `silk_connects`, `silk_backends`, and `silk_routes` functions.

```
shardman.silk_rbc_snap()
```

Retrieves a consistent snapshot of all the connects, backends and routes that can be used by `silk_connects`, `silk_backends`, and `silk_routes` functions.

70.3. Distributed Deadlock Detection

Distributed deadlocks may occur during the processing of distributed transactions. Let us consider the following example:

```
create table players(id int, username text, pass text) with (distributed_by='id');
insert into players select id, 'user_' || id, 'pass_' || id from
generate_series(1,1000) id;
```

Assume that the record with `id=2` belongs to `node1` and the record with `id=3` belongs to `node2`.

Let us execute the following commands on different nodes:

```
node1=# begin;
node1=# update players set pass='someval' where id=3;

node2=# begin;
node2=# update players set pass='someval' where id=2;

-- it should stuck because transaction on node1 locked record with id=3
node2=# update players set pass='someval2' where id=3;

-- it should stuck because transaction on node2 locked record with id=2
node1=# update players set pass='someval2' where id=2;
```

A distributed deadlock situation arises when transactions are mutually locked by each other. PostgreSQL has an internal mechanism for deadlock detection, which detects mutual locking between child processes of a single PostgreSQL instance (backend) and resolves it. However, this mechanism is not applicable to the discovered situation because mutual locking is distributed, i.e., backends from different nodes are involved. From the point of view of the PostgreSQL lock manager, there is no deadlock condition because processes of the single instance are not locking each other. Therefore, Postgres Pro Shardman has its own mechanism for distributed deadlock resolution.

We can represent the interaction between processes in the entire cluster as a graph. A graph vertex represents a process (backend), which we can identify with a couple of attributes `{rgid; vxid}`, where `rgid` is the replication group ID, and `vxid` is the virtual transaction ID of the currently executed transaction. Graph edges represent directional connections between vertices. Each connection is directed from the locked process to the locking process.

It is obvious that any process can be locked by only one process. In other words, if the backend is waiting for a lock, it can only wait for a specific lock. On the other hand, a locking process can acquire multiple locks, meaning that it can lock multiple backends simultaneously.

With that said, the lock graph acts as a singly linked list. If this list contains a closed loop, then here is a deadlock condition. To detect a deadlock, it is necessary to build such a list and detect closed loops in it.

The distributed deadlock detector in Postgres Pro Shardman is implemented as a separate task inside the Postgres Pro Shardman monitor. If a process is unable to acquire a lock within a specified amount of time (which is one second by default, but can be adjusted using the `deadlock_timeout` configuration parameter), the internal PostgreSQL deadlock detector attempts to detect a local deadlock. If no local deadlock is found, the distributed deadlock detector is activated.

The distributed deadlock detector builds a graph (list) of locks in the cluster. It queries views `pg_locks`, `pg_prepared_xacts` and `pg_stat_activity` on the local node and on each of the remote cluster nodes.

The process of building the list of locks involves sequentially querying nodes in the cluster, and it is not atomic, so the list is not consistent. This means that the distributed deadlock detector may produce false positives. During the building of the list, we can store a lock that can disappear before the end of the list building process. To guarantee the reliability of deadlock detection, after the detection of a closed loop, it is necessary to re-query the nodes involved in the closed loop.

After finding the closed loop, the distributed deadlock detector chooses the process belonging to the local node and cancels it. The user process served by the cancelled backend will receive a message:

```
canceled statement due distributed deadlock was found
```

A verbose message about the detected deadlock will be recorded in the server logs:

```
LOG:  distributed deadlock detected
DETAIL:  repgroup 1, PID 95264 (application 'psql'), executed query 'update
players set pass='qqq' where id=2;' is blocked by repgroup 1, PID 95283 (application
'pgfdw:2:95278:9/2'), executed query 'UPDATE public.players_0 SET pass = 'qqq'::text
WHERE ((id = 2))'
    repgroup 1, PID 95283 (application 'pgfdw:2:95278:9/2'), executed query 'UPDATE
public.players_0 SET pass = 'qqq'::text WHERE ((id = 2))' is blocked by repgroup 2,
PID 95278 (application 'psql'), executed query 'update players set pass='qqq' where
id=3;'
    repgroup 2, PID 95278 (application 'psql'), executed query 'update players
set pass='qqq' where id=3;' is blocked by repgroup 2, PID 95267 (application
'pgfdw:1:95264:8/4'), executed query 'UPDATE public.players_1 SET pass = 'qqq'::text
WHERE ((id = 3))'
    repgroup 2, PID 95267 (application 'pgfdw:1:95264:8/4'), executed query 'UPDATE
public.players_1 SET pass = 'qqq'::text WHERE ((id = 3))' is blocked by repgroup 1,
PID 95264 (application 'psql'), executed query 'update players set pass='qqq' where
id=2;'
```

70.4. Global Sequences

Global sequences in Postgres Pro Shardman are implemented on top of regular PostgreSQL sequences with some additional cluster-wide metadata, which among other things holds the interval of globally unused sequence elements.

When `CREATE SEQUENCE .. WITH (global)` is issued, an ordinary PostgreSQL sequence with the same name is created on every cluster node. The range of this local sequence is a bounded sub-interval of the global sequence (as defined by `MINVALUE` and `MAXVALUE` parameters), and it contains at most `block_size` elements. The `nextval` function returns values from the local sequence until it runs out, then a new sub-interval with `block_size` elements is allocated from the global sequence using a broadcast query involving all cluster nodes. So, smaller block size values make the generated numbers more monotonic across the cluster, but incur a performance penalty since the broadcast query may be rather expensive.

Another way to describe the block size parameter is to say that it controls the size of the second cache level, similarly to how the `CACHE` parameter works, except at the level of an entire Postgres Pro Shardman cluster.

Also note, that every time a new sub-interval is allocated the underlying local sequence is modified (as in `ALTER SEQUENCE`), which will lock it for the transaction duration, preventing any other local concurrent transactions from obtaining next sequence values.

70.5. Syncpoints and Consistent Backup

To ensure that cluster binary backup is consistent, Postgres Pro Shardman implements the *syncpoints* mechanism.

To achieve consistent visibility of distributed transactions, the technique of global snapshots based on physical clocks is used. Similarly, it is possible to get a consistent snapshot for backups, only the time corresponding to the global snapshot must be mapped to a set of LSN for each node. Such a set of consistent LSN in a cluster is called a *syncpoint*.

In a Postgres Pro Shardman cluster, each node can generate its own independent local CSN, which does not guarantee the global ordering of values in time. Therefore, we cannot take this arbitrary local CSN as the basis for a *syncpoint*. Instead, Postgres Pro Shardman chooses only those CSNs that match distributed transaction commit records as the basis of the syncpoint. These CSNs have the property of global ordering and can be used to obtain a syncpoint. The main points of this mechanism are described below.

The commit record of each completed transaction in Postgres Pro Shardman contains the assigned CSN for this transaction. This value, together with the LSN of this record, forms a pair of values (CSN, LSN). Each of the cluster nodes stores a certain number of such pairs in RAM in a special structure - the CSNLSNMap. CSNLSNMap is a circular buffer. Each element of the map is a (CSN, LSN) pair. The map size is set by the configuration settings `csn_lsn_map_size`. A (CSN, LSN) pair can be added to the map only if there are no transactions on the node that can receive a CSN less than the one added. This important condition guarantees monotonous growth of CSN and LSN in CSNLSNmap, but does not guarantee that every commit record will get into the map.

When a user submits a request to create a *syncpoint*, a search by every CSNLSNMap is made for a largest possible CSN_g for which there is an entry (CSN_n, LSN) in each node and the condition $CSN_n \leq CSN_g$ is true. The monotonic growth property of every CSNLSNMap ensures that each found pair (CSN_n, LSN) corresponds to the state of the global data at the time corresponding to CSN_g . If no such value of CSN_g is found, the *get syncpoint* operation fails and can be retried later. If such a value CSN_g is found, then a syncpoint is generated as a special type of WAL record, which is duplicated on all nodes of the cluster.

By getting a *syncpoint* and taking the LSN for each node in the cluster from it, we can make a backup of each node, which must necessarily contain that LSN. We can also recover to this LSN using the point in time recovery (PITR) mechanism.

70.6. Collecting Distributed Statement Statistics Using the pgpro_stats Extension

During execution of distributed queries, Postgres Pro Shardman sends derived SQL queries to remote nodes that hold data partitions involved in the query execution. Let's call these SQL queries query fragments. Postgres Pro Shardman sends such queries using the `postgres_fdw` extension. The node that queries the sharded table is called the *coordinator*, while the nodes that accept query fragments are called *shards*. You can use the `shardman.am_coordinator()` function to check whether the current node is the query coordinator.

When the `pgpro_stats` extension is enabled on a Postgres Pro Shardman cluster node, it collects statistics about local and distributed queries. The information about distributed queries initiated by this node is incomplete because it misses data about remote query fragments. The statistics concerning queries

initiated by other nodes is also ambiguous because there is no simple way for a user to determine the distributed query to which the fragment corresponds.

To address these issues, `pgpro_stats` for Postgres Pro Shardman introduces an aggregation of statistics for the distributed queries. These aggregated statistics can be accessed with the `pgpro_stats_sdm_statements` view. However, each Postgres Pro Shardman node collects statistics for all the statements, so that the `pgpro_stats_statements` view can work the way it did before.

When a node receives a query fragment, it saves its statistics to a separate shared hash table. Periodically and asynchronously, each node sends this information from a separate table to the coordinator corresponding to the query. The coordinator aggregates the statistical data obtained from the query fragments with the statistics of its parent query, which is the query initiated by the client.

The `pgpro_stats` extension starts a separate background worker. This worker is responsible for sending the accumulated statistics to the coordinator nodes either every 5 seconds or when triggered by the guard latch. The collecting function sets this latch when the hash table is almost full.

To reduce the network traffic initiated by a statistics sender, compression is applied to the statistics data sent. The compression method can be selected by the `pgpro_stats.transport_compression` configuration parameter.

Each node stores the total number of statistics entries received from the shard node and the timestamp of when they were last received. When a coordinator node receives a statistics message, it updates the appropriate values, which are accessible using the SQL interface.

70.7. Advisory Locks

Postgres Pro Shardman provides API for application to get cluster-wide advisory locks. They are similar to PostgreSQL advisory locks, but are respected on all cluster nodes. Advisory locks and global cluster-wide advisory locks work simultaneously and do not conflict with each other, as global cluster locks use its own locktag type. Both these lock types are shown in the `pg_locks` view and have the `shardman` value in `locktype`.

When global lock is taken, a coordinator sends a lock request to each node in a predefined order (based on `rgid`). When it's a usual lock request, the lock is considered granted when the statement succeeds on all nodes. When it's a try-lock request, the lock is also considered granted when it could be taken on all nodes. If one of the nodes couldn't take the lock, the try-lock function returns `false`. If one of the nodes returned result is different from earlier received (`false` after `true` or `true` after `false`), it is considered to be an error.

To see cluster-wide advisory lock functions, refer to [Advisory Lock Functions for Distributed System](#).

Part VIII. Appendixes

Appendix A. Postgres Pro Error Codes

All messages emitted by the Postgres Pro server are assigned five-character error codes that follow the SQL standard's conventions for "SQLSTATE" codes. Applications that need to know which error condition has occurred should usually test the error code, rather than looking at the textual error message. The error codes are less likely to change across Postgres Pro releases, and also are not subject to change due to localization of error messages. Note that some, but not all, of the error codes produced by Postgres Pro are defined by the SQL standard; some additional error codes for conditions not defined by the standard have been invented or borrowed from other databases.

According to the standard, the first two characters of an error code denote a class of errors, while the last three characters indicate a specific condition within that class. Thus, an application that does not recognize the specific error code might still be able to infer what to do from the error class.

[Table A.1](#) lists all the error codes defined in Postgres Pro Shardman 17.5.1. (Some are not actually used at present, but are defined by the SQL standard.) The error classes are also shown. For each error class there is a "standard" error code having the last three characters 000. This code is used only for error conditions that fall within the class but do not have any more-specific code assigned.

The symbol shown in the column "Condition Name" is the condition name to use in PL/pgSQL. Condition names can be written in either upper or lower case. (Note that PL/pgSQL does not recognize warning, as opposed to error, condition names; those are classes 00, 01, and 02.)

For some types of errors, the server reports the name of a database object (a table, table column, data type, or constraint) associated with the error; for example, the name of the unique constraint that caused a `unique_violation` error. Such names are supplied in separate fields of the error report message so that applications need not try to extract them from the possibly-localized human-readable text of the message. As of PostgreSQL 9.3, complete coverage for this feature exists only for errors in SQLSTATE class 23 (integrity constraint violation), but this is likely to be expanded in future.

Table A.1. Postgres Pro Error Codes

Error Code	Condition Name
Class 00 — Successful Completion	
00000	successful_completion
Class 01 — Warning	
01000	warning
0100C	dynamic_result_sets_returned
01008	implicit_zero_bit_padding
01003	null_value_eliminated_in_set_function
01007	privilege_not_granted
01006	privilege_not_revoked
01004	string_data_right_truncation
01P01	deprecated_feature
Class 02 — No Data (this is also a warning class per the SQL standard)	
02000	no_data
02001	no_additional_dynamic_result_sets_returned
Class 03 — SQL Statement Not Yet Complete	
03000	sql_statement_not_yet_complete
Class 08 — Connection Exception	
08000	connection_exception

Error Code	Condition Name
08003	connection_does_not_exist
08006	connection_failure
08001	sqlclient_unable_to_establish_sqlconnection
08004	sqlserver_rejected_establishment_of_sqlconnection
08007	transaction_resolution_unknown
08P01	protocol_violation
Class 09 — Triggered Action Exception	
09000	triggered_action_exception
Class 0A — Feature Not Supported	
0A000	feature_not_supported
Class 0B — Invalid Transaction Initiation	
0B000	invalid_transaction_initiation
Class 0F — Locator Exception	
0F000	locator_exception
0F001	invalid_locator_specification
Class 0L — Invalid Grantor	
0L000	invalid_grantor
0LP01	invalid_grant_operation
Class 0P — Invalid Role Specification	
0P000	invalid_role_specification
Class 0Z — Diagnostics Exception	
0Z000	diagnostics_exception
0Z002	stacked_diagnostics_accessed_without_active_handler
Class 20 — Case Not Found	
20000	case_not_found
Class 21 — Cardinality Violation	
21000	cardinality_violation
Class 22 — Data Exception	
22000	data_exception
2202E	array_subscript_error
22021	character_not_in_repertoire
22008	datetime_field_overflow
22012	division_by_zero
22005	error_in_assignment
2200B	escape_character_conflict
22022	indicator_overflow
22015	interval_field_overflow
2201E	invalid_argument_for_logarithm
22014	invalid_argument_for_ntile_function

Postgres Pro Error Codes

Error Code	Condition Name
22016	invalid_argument_for_nth_value_function
2201F	invalid_argument_for_power_function
2201G	invalid_argument_for_width_bucket_function
22018	invalid_character_value_for_cast
22007	invalid_datetime_format
22019	invalid_escape_character
2200D	invalid_escape_octet
22025	invalid_escape_sequence
22P06	nonstandard_use_of_escape_character
22010	invalid_indicator_parameter_value
22023	invalid_parameter_value
22013	invalid_preceding_or_following_size
2201B	invalid_regular_expression
2201W	invalid_row_count_in_limit_clause
2201X	invalid_row_count_in_result_offset_clause
2202H	invalid_tablesample_argument
2202G	invalid_tablesample_repeat
22009	invalid_time_zone_displacement_value
2200C	invalid_use_of_escape_character
2200G	most_specific_type_mismatch
22004	null_value_not_allowed
22002	null_value_no_indicator_parameter
22003	numeric_value_out_of_range
2200H	sequence_generator_limit_exceeded
22026	string_data_length_mismatch
22001	string_data_right_truncation
22011	substring_error
22027	trim_error
22024	unterminated_c_string
2200F	zero_length_character_string
22P01	floating_point_exception
22P02	invalid_text_representation
22P03	invalid_binary_representation
22P04	bad_copy_file_format
22P05	untranslatable_character
2200L	not_an_xml_document
2200M	invalid_xml_document
2200N	invalid_xml_content
2200S	invalid_xml_comment

Error Code	Condition Name
2200T	invalid_xml_processing_instruction
22030	duplicate_json_object_key_value
22031	invalid_argument_for_sql_json_datetime_function
22032	invalid_json_text
22033	invalid_sql_json_subscript
22034	more_than_one_sql_json_item
22035	no_sql_json_item
22036	non_numeric_sql_json_item
22037	non_unique_keys_in_a_json_object
22038	singleton_sql_json_item_required
22039	sql_json_array_not_found
2203A	sql_json_member_not_found
2203B	sql_json_number_not_found
2203C	sql_json_object_not_found
2203D	too_many_json_array_elements
2203E	too_many_json_object_members
2203F	sql_json_scalar_required
2203G	sql_json_item_cannot_be_cast_to_target_type
Class 23 — Integrity Constraint Violation	
23000	integrity_constraint_violation
23001	restrict_violation
23502	not_null_violation
23503	foreign_key_violation
23505	unique_violation
23514	check_violation
23P01	exclusion_violation
Class 24 — Invalid Cursor State	
24000	invalid_cursor_state
Class 25 — Invalid Transaction State	
25000	invalid_transaction_state
25001	active_sql_transaction
25002	branch_transaction_already_active
25008	held_cursor_requires_same_isolation_level
25003	inappropriate_access_mode_for_branch_transaction
25004	inappropriate_isolation_level_for_branch_transaction
25005	no_active_sql_transaction_for_branch_transaction
25006	read_only_sql_transaction
25007	schema_and_data_statement_mixing_not_supported
25P01	no_active_sql_transaction

Error Code	Condition Name
25P02	in_failed_sql_transaction
25P03	idle_in_transaction_session_timeout
25P04	transaction_timeout
25P99	query_inadequate_execution_time
Class 26 – Invalid SQL Statement Name	
26000	invalid_sql_statement_name
Class 27 – Triggered Data Change Violation	
27000	triggered_data_change_violation
Class 28 – Invalid Authorization Specification	
28000	invalid_authorization_specification
28P01	invalid_password
Class 2B – Dependent Privilege Descriptors Still Exist	
2B000	dependent_privilege_descriptors_still_exist
2BP01	dependent_objects_still_exist
Class 2D – Invalid Transaction Termination	
2D000	invalid_transaction_termination
Class 2F – SQL Routine Exception	
2F000	sql_routine_exception
2F005	function_executed_no_return_statement
2F002	modifying_sql_data_not_permitted
2F003	prohibited_sql_statement_attempted
2F004	reading_sql_data_not_permitted
Class 34 – Invalid Cursor Name	
34000	invalid_cursor_name
Class 38 – External Routine Exception	
38000	external_routine_exception
38001	containing_sql_not_permitted
38002	modifying_sql_data_not_permitted
38003	prohibited_sql_statement_attempted
38004	reading_sql_data_not_permitted
Class 39 – External Routine Invocation Exception	
39000	external_routine_invocation_exception
39001	invalid_sqlstate_returned
39004	null_value_not_allowed
39P01	trigger_protocol_violated
39P02	srf_protocol_violated
39P03	event_trigger_protocol_violated
Class 3B – Savepoint Exception	
3B000	savepoint_exception

Error Code	Condition Name
3B001	invalid_savepoint_specification
Class 3D — Invalid Catalog Name	
3D000	invalid_catalog_name
Class 3F — Invalid Schema Name	
3F000	invalid_schema_name
Class 40 — Transaction Rollback	
40000	transaction_rollback
40002	transaction_integrity_constraint_violation
40001	serialization_failure
40003	statement_completion_unknown
40P01	deadlock_detected
Class 42 — Syntax Error or Access Rule Violation	
42000	syntax_error_or_access_rule_violation
42601	syntax_error
42501	insufficient_privilege
42846	cannot_coerce
42803	grouping_error
42P20	windowing_error
42P19	invalid_recursion
42830	invalid_foreign_key
42602	invalid_name
42622	name_too_long
42939	reserved_name
42804	datatype_mismatch
42P18	indeterminate_datatype
42P21	collation_mismatch
42P22	indeterminate_collation
42809	wrong_object_type
428C9	generated_always
42703	undefined_column
42883	undefined_function
42P01	undefined_table
42P02	undefined_parameter
42704	undefined_object
42701	duplicate_column
42P03	duplicate_cursor
42P04	duplicate_database
42723	duplicate_function
42P05	duplicate_prepared_statement

Error Code	Condition Name
42P06	duplicate_schema
42P07	duplicate_table
42712	duplicate_alias
42710	duplicate_object
42702	ambiguous_column
42725	ambiguous_function
42P08	ambiguous_parameter
42P09	ambiguous_alias
42P10	invalid_column_reference
42611	invalid_column_definition
42P11	invalid_cursor_definition
42P12	invalid_database_definition
42P13	invalid_function_definition
42P14	invalid_prepared_statement_definition
42P15	invalid_schema_definition
42P16	invalid_table_definition
42P17	invalid_object_definition
Class 44 — WITH CHECK OPTION Violation	
44000	with_check_option_violation
Class 53 — Insufficient Resources	
53000	insufficient_resources
53100	disk_full
53200	out_of_memory
53300	too_many_connections
53400	configuration_limit_exceeded
53500	cardinality_estimation_accuracy
Class 54 — Program Limit Exceeded	
54000	program_limit_exceeded
54001	statement_too_complex
54011	too_many_columns
54023	too_many_arguments
Class 55 — Object Not In Prerequisite State	
55000	object_not_in_prerequisite_state
55006	object_in_use
55P02	cant_change_runtime_param
55P03	lock_not_available
55P04	unsafe_new_enum_value_usage
Class 57 — Operator Intervention	
57000	operator_intervention

Error Code	Condition Name
57014	query_canceled
57P01	admin_shutdown
57P02	crash_shutdown
57P03	cannot_connect_now
57P04	database_dropped
57P05	idle_session_timeout
Class 58 — System Error (errors external to PostgreSQL itself)	
58000	system_error
58030	io_error
58P01	undefined_file
58P02	duplicate_file
Class F0 — Configuration File Error	
F0000	config_file_error
F0001	lock_file_exists
Class HV — Foreign Data Wrapper Error (SQL/MED)	
HV000	fdw_error
HV005	fdw_column_name_not_found
HV002	fdw_dynamic_parameter_value_needed
HV010	fdw_function_sequence_error
HV021	fdw_inconsistent_descriptor_information
HV024	fdw_invalid_attribute_value
HV007	fdw_invalid_column_name
HV008	fdw_invalid_column_number
HV004	fdw_invalid_data_type
HV006	fdw_invalid_data_type_descriptors
HV091	fdw_invalid_descriptor_field_identifier
HV00B	fdw_invalid_handle
HV00C	fdw_invalid_option_index
HV00D	fdw_invalid_option_name
HV090	fdw_invalid_string_length_or_buffer_length
HV00A	fdw_invalid_string_format
HV009	fdw_invalid_use_of_null_pointer
HV014	fdw_too_many_handles
HV001	fdw_out_of_memory
HV00P	fdw_no_schemas
HV00J	fdw_option_name_not_found
HV00K	fdw_reply_handle
HV00Q	fdw_schema_not_found
HV00R	fdw_table_not_found

Error Code	Condition Name
HV00L	fdw_unable_to_create_execution
HV00M	fdw_unable_to_create_reply
HV00N	fdw_unable_to_establish_connection
Class P0 — PL/pgSQL Error	
P0000	plpgsql_error
P0001	raise_exception
P0002	no_data_found
P0003	too_many_rows
P0004	assert_failure
Class XX — Internal Error	
XX000	internal_error
XX001	data_corrupted
XX002	index_corrupted

Appendix B. Date/Time Support

Postgres Pro uses an internal heuristic parser for all date/time input support. Dates and times are input as strings, and are broken up into distinct fields with a preliminary determination of what kind of information can be in the field. Each field is interpreted and either assigned a numeric value, ignored, or rejected. The parser contains internal lookup tables for all textual fields, including months, days of the week, and time zones.

This appendix includes information on the content of these lookup tables and describes the steps used by the parser to decode dates and times.

B.1. Date/Time Input Interpretation

Date/time input strings are decoded using the following procedure.

1. Break the input string into tokens and categorize each token as a string, time, time zone, or number.
 - a. If the numeric token contains a colon (:), this is a time string. Include all subsequent digits and colons.
 - b. If the numeric token contains a dash (-), slash (/), or two or more dots (.), this is a date string which might have a text month. If a date token has already been seen, it is instead interpreted as a time zone name (e.g., `America/New_York`).
 - c. If the token is numeric only, then it is either a single field or an ISO 8601 concatenated date (e.g., `19990113` for January 13, 1999) or time (e.g., `141516` for 14:15:16).
 - d. If the token starts with a plus (+) or minus (-), then it is either a numeric time zone or a special field.
2. If the token is an alphabetic string, match up with possible strings:
 - a. See if the token matches any known time zone abbreviation. These abbreviations are supplied by the configuration file described in [Section B.4](#).
 - b. If not found, search an internal table to match the token as either a special string (e.g., `today`), day (e.g., `Thursday`), month (e.g., `January`), or noise word (e.g., `at`, `on`).
 - c. If still not found, throw an error.
3. When the token is a number or number field:
 - a. If there are eight or six digits, and if no other date fields have been previously read, then interpret as a “concatenated date” (e.g., `19990118` or `990118`). The interpretation is `YYYYMMDD` or `YYMMDD`.
 - b. If the token is three digits and a year has already been read, then interpret as day of year.
 - c. If four or six digits and a year has already been read, then interpret as a time (`HHMM` or `HHMMSS`).
 - d. If three or more digits and no date fields have yet been found, interpret as a year (this forces `yy-mm-dd` ordering of the remaining date fields).
 - e. Otherwise the date field ordering is assumed to follow the `DateStyle` setting: `mm-dd-yy`, `dd-mm-yy`, or `yy-mm-dd`. Throw an error if a month or day field is found to be out of range.
4. If BC has been specified, negate the year and add one for internal storage. (There is no year zero in the Gregorian calendar, so numerically 1 BC becomes year zero.)
5. If BC was not specified, and if the year field was two digits in length, then adjust the year to four digits. If the field is less than 70, then add 2000, otherwise add 1900.

Tip

Gregorian years AD 1-99 can be entered by using 4 digits with leading zeros (e.g., 0099 is AD 99).

B.2. Handling of Invalid or Ambiguous Timestamps

Ordinarily, if a date/time string is syntactically valid but contains out-of-range field values, an error will be thrown. For example, input specifying the 31st of February will be rejected.

During a daylight-savings-time transition, it is possible for a seemingly valid timestamp string to represent a nonexistent or ambiguous timestamp. Such cases are not rejected; the ambiguity is resolved by determining which UTC offset to apply. For example, supposing that the [TimeZone](#) parameter is set to `America/New_York`, consider

```
=> SELECT '2018-03-11 02:30'::timestampz;
       timestampz
-----
2018-03-11 03:30:00-04
(1 row)
```

Because that day was a spring-forward transition date in that time zone, there was no civil time instant 2:30AM; clocks jumped forward from 2AM EST to 3AM EDT. Postgres Pro interprets the given time as if it were standard time (UTC-5), which then renders as 3:30AM EDT (UTC-4).

Conversely, consider the behavior during a fall-back transition:

```
=> SELECT '2018-11-04 01:30'::timestampz;
       timestampz
-----
2018-11-04 01:30:00-05
(1 row)
```

On that date, there were two possible interpretations of 1:30AM; there was 1:30AM EDT, and then an hour later after clocks jumped back from 2AM EDT to 1AM EST, there was 1:30AM EST. Again, Postgres Pro interprets the given time as if it were standard time (UTC-5). We can force the other interpretation by specifying daylight-savings time:

```
=> SELECT '2018-11-04 01:30 EDT'::timestampz;
       timestampz
-----
2018-11-04 01:30:00-04
(1 row)
```

The precise rule that is applied in such cases is that an invalid timestamp that appears to fall within a jump-forward daylight savings transition is assigned the UTC offset that prevailed in the time zone just before the transition, while an ambiguous timestamp that could fall on either side of a jump-back transition is assigned the UTC offset that prevailed just after the transition. In most time zones this is equivalent to saying that “the standard-time interpretation is preferred when in doubt”.

In all cases, the UTC offset associated with a timestamp can be specified explicitly, using either a numeric UTC offset or a time zone abbreviation that corresponds to a fixed UTC offset. The rule just given applies only when it is necessary to infer a UTC offset for a time zone in which the offset varies.

B.3. Date/Time Key Words

[Table B.1](#) shows the tokens that are recognized as names of months.

Table B.1. Month Names

Month	Abbreviations
January	Jan
February	Feb
March	Mar
April	Apr
May	
June	Jun
July	Jul
August	Aug
September	Sep, Sept
October	Oct
November	Nov
December	Dec

Table B.2 shows the tokens that are recognized as names of days of the week.

Table B.2. Day of the Week Names

Day	Abbreviations
Sunday	Sun
Monday	Mon
Tuesday	Tue, Tues
Wednesday	Wed, Weds
Thursday	Thu, Thur, Thurs
Friday	Fri
Saturday	Sat

Table B.3 shows the tokens that serve various modifier purposes.

Table B.3. Date/Time Field Modifiers

Identifier	Description
AM	Time is before 12:00
AT	Ignored
JULIAN, JD, J	Next field is Julian Date
ON	Ignored
PM	Time is on or after 12:00
T	Next field is time

B.4. Date/Time Configuration Files

Since timezone abbreviations are not well standardized, Postgres Pro provides a means to customize the set of abbreviations accepted by the server. The [timezone_abbreviations](#) run-time parameter determines the active set of abbreviations. While this parameter can be altered by any database user, the possible values for it are under the control of the database administrator — they are in fact names of configuration files stored in `.../share/timezonesets/` of the installation directory. By adding or altering files in that directory, the administrator can set local policy for timezone abbreviations.

`timezone_abbreviations` can be set to any file name found in `.../share/zoneinfo/`, if the file's name is entirely alphabetic. (The prohibition against non-alphabetic characters in `timezone_abbreviations` prevents reading files outside the intended directory, as well as reading editor backup files and other extraneous files.)

A timezone abbreviation file can contain blank lines and comments beginning with `#`. Non-comment lines must have one of these formats:

```
zone_abbreviation offset
zone_abbreviation offset D
zone_abbreviation time_zone_name
@INCLUDE file_name
@OVERRIDE
```

A *zone_abbreviation* is just the abbreviation being defined. An *offset* is an integer giving the equivalent offset in seconds from UTC, positive being east from Greenwich and negative being west. For example, -18000 would be five hours west of Greenwich, or North American east coast standard time. `D` indicates that the zone name represents local daylight-savings time rather than standard time.

Alternatively, a *time_zone_name* can be given, referencing a zone name defined in the IANA timezone database. The zone's definition is consulted to see whether the abbreviation is or has been in use in that zone, and if so, the appropriate meaning is used — that is, the meaning that was currently in use at the timestamp whose value is being determined, or the meaning in use immediately before that if it wasn't current at that time, or the oldest meaning if it was used only after that time. This behavior is essential for dealing with abbreviations whose meaning has historically varied. It is also allowed to define an abbreviation in terms of a zone name in which that abbreviation does not appear; then using the abbreviation is just equivalent to writing out the zone name.

Tip

Using a simple integer *offset* is preferred when defining an abbreviation whose offset from UTC has never changed, as such abbreviations are much cheaper to process than those that require consulting a time zone definition.

The `@INCLUDE` syntax allows inclusion of another file in the `.../share/zoneinfo/` directory. Inclusion can be nested, to a limited depth.

The `@OVERRIDE` syntax indicates that subsequent entries in the file can override previous entries (typically, entries obtained from included files). Without this, conflicting definitions of the same timezone abbreviation are considered an error.

In an unmodified installation, the file `Default` contains all the non-conflicting time zone abbreviations for most of the world. Additional files `Australia` and `India` are provided for those regions: these files first include the `Default` file and then add or modify abbreviations as needed.

For reference purposes, a standard installation also contains files `Africa.txt`, `America.txt`, etc., containing information about every time zone abbreviation known to be in use according to the IANA timezone database. The zone name definitions found in these files can be copied and pasted into a custom configuration file as needed. Note that these files cannot be directly referenced as `timezone_abbreviations` settings, because of the dot embedded in their names.

Note

If an error occurs while reading the time zone abbreviation set, no new value is applied and the old set is kept. If the error occurs while starting the database, startup fails.

Caution

Time zone abbreviations defined in the configuration file override non-timezone meanings built into Postgres Pro. For example, the `Australia` configuration file defines `SAT` (for South Australian Standard Time). When this file is active, `SAT` will not be recognized as an abbreviation for Saturday.

Caution

If you modify files in `.../share/timezonesets/`, it is up to you to make backups — a normal database dump will not include this directory.

B.5. POSIX Time Zone Specifications

Postgres Pro can accept time zone specifications that are written according to the POSIX standard's rules for the `TZ` environment variable. POSIX time zone specifications are inadequate to deal with the complexity of real-world time zone history, but there are sometimes reasons to use them.

A POSIX time zone specification has the form

```
STD offset [ DST [ dstoffset ] [ , rule ] ]
```

(For readability, we show spaces between the fields, but spaces should not be used in practice.) The fields are:

- *STD* is the zone abbreviation to be used for standard time.
- *offset* is the zone's standard-time offset from UTC.
- *DST* is the zone abbreviation to be used for daylight-savings time. If this field and the following ones are omitted, the zone uses a fixed UTC offset with no daylight-savings rule.
- *dstoffset* is the daylight-savings offset from UTC. This field is typically omitted, since it defaults to one hour less than the standard-time *offset*, which is usually the right thing.
- *rule* defines the rule for when daylight savings is in effect, as described below.

In this syntax, a zone abbreviation can be a string of letters, such as `EST`, or an arbitrary string surrounded by angle brackets, such as `<UTC-05>`. Note that the zone abbreviations given here are only used for output, and even then only in some timestamp output formats. The zone abbreviations recognized in timestamp input are determined as explained in [Section B.4](#).

The offset fields specify the hours, and optionally minutes and seconds, difference from UTC. They have the format `hh[:mm[:ss]]` optionally with a leading sign (+ or -). The positive sign is used for zones west of Greenwich. (Note that this is the opposite of the ISO-8601 sign convention used elsewhere in Postgres Pro.) *hh* can have one or two digits; *mm* and *ss* (if used) must have two.

The daylight-savings transition *rule* has the format

```
dstdate [ / dsttime ] , stddate [ / stdtime ]
```

(As before, spaces should not be included in practice.) The *dstdate* and *dsttime* fields define when daylight-savings time starts, while *stddate* and *stdtime* define when standard time starts. (In some cases, notably in zones south of the equator, the former might be later in the year than the latter.) The date fields have one of these formats:

n

A plain integer denotes a day of the year, counting from zero to 364, or to 365 in leap years.

Jn

In this form, *n* counts from 1 to 365, and February 29 is not counted even if it is present. (Thus, a transition occurring on February 29 could not be specified this way. However, days after February

have the same numbers whether it's a leap year or not, so that this form is usually more useful than the plain-integer form for transitions on fixed dates.)

Mm.n.d

This form specifies a transition that always happens during the same month and on the same day of the week. *m* identifies the month, from 1 to 12. *n* specifies the *n*'th occurrence of the weekday identified by *d*. *n* is a number between 1 and 4, or 5 meaning the last occurrence of that weekday in the month (which could be the fourth or the fifth). *d* is a number between 0 and 6, with 0 indicating Sunday. For example, `M3.2.0` means “the second Sunday in March”.

Note

The `M` format is sufficient to describe many common daylight-savings transition laws. But note that none of these variants can deal with daylight-savings law changes, so in practice the historical data stored for named time zones (in the IANA time zone database) is necessary to interpret past time stamps correctly.

The time fields in a transition rule have the same format as the offset fields described previously, except that they cannot contain signs. They define the current local time at which the change to the other time occurs. If omitted, they default to `02:00:00`.

If a daylight-savings abbreviation is given but the transition *rule* field is omitted, the fallback behavior is to use the rule `M3.2.0,M11.1.0`, which corresponds to USA practice as of 2020 (that is, spring forward on the second Sunday of March, fall back on the first Sunday of November, both transitions occurring at 2AM prevailing time). Note that this rule does not give correct USA transition dates for years before 2007.

As an example, `CET-1CEST,M3.5.0,M10.5.0/3` describes the current (as of 2020) timekeeping practice in Paris. This specification says that standard time has the abbreviation `CET` and is one hour ahead (east) of UTC; daylight savings time has the abbreviation `CEST` and is implicitly two hours ahead of UTC; daylight savings time begins on the last Sunday in March at 2AM CET and ends on the last Sunday in October at 3AM CEST.

The four timezone names `EST5EDT`, `CST6CDT`, `MST7MDT`, and `PST8PDT` look like they are POSIX zone specifications. However, they actually are treated as named time zones because (for historical reasons) there are files by those names in the IANA time zone database. The practical implication of this is that these zone names will produce valid historical USA daylight-savings transitions, even when a plain POSIX specification would not.

One should be wary that it is easy to misspell a POSIX-style time zone specification, since there is no check on the reasonableness of the zone abbreviation(s). For example, `SET TIMEZONE TO FOOBAR0` will work, leaving the system effectively using a rather peculiar abbreviation for UTC.

B.6. History of Units

The SQL standard states that “Within the definition of a ‘datetime literal’, the ‘datetime values’ are constrained by the natural rules for dates and times according to the Gregorian calendar”. Postgres Pro follows the SQL standard's lead by counting dates exclusively in the Gregorian calendar, even for years before that calendar was in use. This rule is known as the *proleptic Gregorian calendar*.

The Julian calendar was introduced by Julius Caesar in 45 BC. It was in common use in the Western world until the year 1582, when countries started changing to the Gregorian calendar. In the Julian calendar, the tropical year is approximated as $365 \frac{1}{4}$ days = 365.25 days. This gives an error of about 1 day in 128 years.

The accumulating calendar error prompted Pope Gregory XIII to reform the calendar in accordance with instructions from the Council of Trent. In the Gregorian calendar, the tropical year is approximated as

$365 + 97 / 400$ days = 365.2425 days. Thus it takes approximately 3300 years for the tropical year to shift one day with respect to the Gregorian calendar.

The approximation $365+97/400$ is achieved by having 97 leap years every 400 years, using the following rules:

Every year divisible by 4 is a leap year.

However, every year divisible by 100 is not a leap year.

However, every year divisible by 400 is a leap year after all.

So, 1700, 1800, 1900, 2100, and 2200 are not leap years. But 1600, 2000, and 2400 are leap years. By contrast, in the older Julian calendar all years divisible by 4 are leap years.

The papal bull of February 1582 decreed that 10 days should be dropped from October 1582 so that 15 October should follow immediately after 4 October. This was observed in Italy, Poland, Portugal, and Spain. Other Catholic countries followed shortly after, but Protestant countries were reluctant to change, and the Greek Orthodox countries didn't change until the start of the 20th century. The reform was observed by Great Britain and its dominions (including what is now the USA) in 1752. Thus 2 September 1752 was followed by 14 September 1752. This is why Unix systems that have the `cal` program produce the following:

```
$ cal 9 1752
    September 1752
 S  M Tu   W Th   F   S
      1   2  14  15  16
17 18 19 20 21 22 23
24 25 26 27 28 29 30
```

But, of course, this calendar is only valid for Great Britain and dominions, not other places. Since it would be difficult and confusing to try to track the actual calendars that were in use in various places at various times, Postgres Pro does not try, but rather follows the Gregorian calendar rules for all dates, even though this method is not historically accurate.

Different calendars have been developed in various parts of the world, many predating the Gregorian system. For example, the beginnings of the Chinese calendar can be traced back to the 14th century BC. Legend has it that the Emperor Huangdi invented that calendar in 2637 BC. The People's Republic of China uses the Gregorian calendar for civil purposes. The Chinese calendar is used for determining festivals.

B.7. Julian Dates

The *Julian Date* system is a method for numbering days. It is unrelated to the Julian calendar, though it is confusingly named similarly to that calendar. The Julian Date system was invented by the French scholar Joseph Justus Scaliger (1540-1609) and probably takes its name from Scaliger's father, the Italian scholar Julius Caesar Scaliger (1484-1558).

In the Julian Date system, each day has a sequential number, starting from JD 0 (which is sometimes called *the* Julian Date). JD 0 corresponds to 1 January 4713 BC in the Julian calendar, or 24 November 4714 BC in the Gregorian calendar. Julian Date counting is most often used by astronomers for labeling their nightly observations, and therefore a date runs from noon UTC to the next noon UTC, rather than from midnight to midnight: JD 0 designates the 24 hours from noon UTC on 24 November 4714 BC to noon UTC on 25 November 4714 BC.

Although Postgres Pro supports Julian Date notation for input and output of dates (and also uses Julian dates for some internal datetime calculations), it does not observe the nicety of having dates run from noon to noon. Postgres Pro treats a Julian Date as running from local midnight to local midnight, the same as a normal date.

This definition does, however, provide a way to obtain the astronomical definition when you need it: do the arithmetic in time zone `UTC+12`. For example,


```
=> SELECT extract(julian from '2021-06-23 7:00:00-04'::timestampz at time zone 'UTC
+12');
```

```
extract
```

```
-----
2459388.9583333333333333333333333333
```

```
(1 row)
```

```
=> SELECT extract(julian from '2021-06-23 8:00:00-04'::timestampz at time zone 'UTC
+12');
```

```
extract
```

```
-----
2459389.0000000000000000000000000000000000
```

```
(1 row)
```

```
=> SELECT extract(julian from date '2021-06-23');
```

```
extract
```

```
-----
2459389
```

```
(1 row)
```

Appendix C. SQL Key Words

[Table C.1](#) lists all tokens that are key words in the SQL standard and in Postgres Pro Shardman 17.5.1. Background information can be found in [Section 4.1.1](#). (For space reasons, only the latest two versions of the SQL standard, and SQL-92 for historical comparison, are included. The differences between those and the other intermediate standard versions are small.)

SQL distinguishes between *reserved* and *non-reserved* key words. According to the standard, reserved key words are the only real key words; they are never allowed as identifiers. Non-reserved key words only have a special meaning in particular contexts and can be used as identifiers in other contexts. Most non-reserved key words are actually the names of built-in tables and functions specified by SQL. The concept of non-reserved key words essentially only exists to declare that some predefined meaning is attached to a word in some contexts.

In the Postgres Pro parser, life is a bit more complicated. There are several different classes of tokens ranging from those that can never be used as an identifier to those that have absolutely no special status in the parser, but are considered ordinary identifiers. (The latter is usually the case for functions specified by SQL.) Even reserved key words are not completely reserved in Postgres Pro, but can be used as column labels (for example, `SELECT 55 AS CHECK`, even though `CHECK` is a reserved key word).

In [Table C.1](#) in the column for Postgres Pro we classify as “non-reserved” those key words that are explicitly known to the parser but are allowed as column or table names. Some key words that are otherwise non-reserved cannot be used as function or data type names and are marked accordingly. (Most of these words represent built-in functions or data types with special syntax. The function or type is still available but it cannot be redefined by the user.) Labeled “reserved” are those tokens that are not allowed as column or table names. Some reserved key words are allowable as names for functions or data types; this is also shown in the table. If not so marked, a reserved key word is only allowed as a column label. A blank entry in this column means that the word is treated as an ordinary identifier by Postgres Pro.

Furthermore, while most key words can be used as “bare” column labels without writing `AS` before them (as described in [Section 7.3.2](#)), there are a few that require a leading `AS` to avoid ambiguity. These are marked in the table as “requires `AS`”.

As a general rule, if you get spurious parser errors for commands that use any of the listed key words as an identifier, you should try quoting the identifier to see if the problem goes away.

It is important to understand before studying [Table C.1](#) that the fact that a key word is not reserved in Postgres Pro does not mean that the feature related to the word is not implemented. Conversely, the presence of a key word does not indicate the existence of a feature.

Table C.1. SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
A		non-reserved	non-reserved	
ABORT	non-reserved			
ABS		reserved	reserved	
ABSENT	non-reserved	reserved	reserved	
ABSOLUTE	non-reserved	non-reserved	non-reserved	reserved
ACCESS	non-reserved			
ACCORDING		non-reserved	non-reserved	
ACCOUNT	non-reserved			
ACOS		reserved	reserved	
ACTION	non-reserved	non-reserved	non-reserved	reserved
ADA		non-reserved	non-reserved	non-reserved

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
ADD	non-reserved	non-reserved	non-reserved	reserved
ADMIN	non-reserved	non-reserved	non-reserved	
AFTER	non-reserved	non-reserved	non-reserved	
AGGREGATE	non-reserved			
ALL	reserved	reserved	reserved	reserved
ALLOCATE		reserved	reserved	reserved
ALSO	non-reserved			
ALTER	non-reserved	reserved	reserved	reserved
ALWAYS	non-reserved	non-reserved	non-reserved	
ANALYSE	reserved			
ANALYZE	reserved			
AND	reserved	reserved	reserved	reserved
ANY	reserved	reserved	reserved	reserved
ANY_VALUE		reserved		
APPLICATION	non-reserved			
ARE		reserved	reserved	reserved
ARRAY	reserved, requires AS	reserved	reserved	
ARRAY_AGG		reserved	reserved	
ARRAY_MAX_CARDINALITY		reserved	reserved	
AS	reserved, requires AS	reserved	reserved	reserved
ASC	reserved	non-reserved	non-reserved	reserved
ASENSITIVE	non-reserved	reserved	reserved	
ASIN		reserved	reserved	
ASSERTION	non-reserved	non-reserved	non-reserved	reserved
ASSIGNMENT	non-reserved	non-reserved	non-reserved	
ASYMMETRIC	reserved	reserved	reserved	
AT	non-reserved	reserved	reserved	reserved
ATAN		reserved	reserved	
ATOMIC	non-reserved	reserved	reserved	
ATTACH	non-reserved			
ATTRIBUTE	non-reserved	non-reserved	non-reserved	
ATTRIBUTES		non-reserved	non-reserved	
AUTHORIZATION	reserved (can be function or type)	reserved	reserved	reserved
AUTONOMOUS	reserved (can be function or type)			
AVG		reserved	reserved	reserved
BACKWARD	non-reserved			
BASE64		non-reserved	non-reserved	

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
BEFORE	non-reserved	non-reserved	non-reserved	
BEGIN	non-reserved	reserved	reserved	reserved
BEGIN_FRAME		reserved	reserved	
BEGIN_PARTITION		reserved	reserved	
BERNOULLI		non-reserved	non-reserved	
BETWEEN	non-reserved (cannot be function or type)	reserved	reserved	reserved
BIGINT	non-reserved (cannot be function or type)	reserved	reserved	
BINARY	reserved (can be function or type)	reserved	reserved	
BIT	non-reserved (cannot be function or type)			reserved
BIT_LENGTH				reserved
BLOB		reserved	reserved	
BLOCKED		non-reserved	non-reserved	
BOM		non-reserved	non-reserved	
BOOLEAN	non-reserved (cannot be function or type)	reserved	reserved	
BOTH	reserved	reserved	reserved	reserved
BREADTH	non-reserved	non-reserved	non-reserved	
BTRIM		reserved		
BY	non-reserved	reserved	reserved	reserved
C		non-reserved	non-reserved	non-reserved
CACHE	non-reserved			
CALL	non-reserved	reserved	reserved	
CALLED	non-reserved	reserved	reserved	
CARDINALITY		reserved	reserved	
CASCADE	non-reserved	non-reserved	non-reserved	reserved
CASCADED	non-reserved	reserved	reserved	reserved
CASE	reserved	reserved	reserved	reserved
CAST	reserved	reserved	reserved	reserved
CATALOG	non-reserved	non-reserved	non-reserved	reserved
CATALOG_NAME		non-reserved	non-reserved	non-reserved
CEIL		reserved	reserved	
CEILING		reserved	reserved	
CHAIN	non-reserved	non-reserved	non-reserved	
CHAINING		non-reserved	non-reserved	

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
CHAR	non-reserved (cannot be function or type), requires AS	reserved	reserved	reserved
CHARACTER	non-reserved (cannot be function or type), requires AS	reserved	reserved	reserved
CHARACTERISTICS	non-reserved	non-reserved	non-reserved	
CHARACTERS		non-reserved	non-reserved	
CHARACTER_LENGTH		reserved	reserved	reserved
CHARACTER_SET_CATALOG		non-reserved	non-reserved	non-reserved
CHARACTER_SET_NAME		non-reserved	non-reserved	non-reserved
CHARACTER_SET_SCHEMA		non-reserved	non-reserved	non-reserved
CHAR_LENGTH		reserved	reserved	reserved
CHECK	reserved	reserved	reserved	reserved
CHECKPOINT	non-reserved			
CLASS	non-reserved			
CLASSIFIER		reserved	reserved	
CLASS_ORIGIN		non-reserved	non-reserved	non-reserved
CLOB		reserved	reserved	
CLOSE	non-reserved	reserved	reserved	reserved
CLUSTER	non-reserved			
COALESCE	non-reserved (cannot be function or type)	reserved	reserved	reserved
COBOL		non-reserved	non-reserved	non-reserved
COLLATE	reserved	reserved	reserved	reserved
COLLATION	reserved (can be function or type)	non-reserved	non-reserved	reserved
COLLATION_CATALOG		non-reserved	non-reserved	non-reserved
COLLATION_NAME		non-reserved	non-reserved	non-reserved
COLLATION_SCHEMA		non-reserved	non-reserved	non-reserved
COLLECT		reserved	reserved	
COLUMN	reserved	reserved	reserved	reserved
COLUMNS	non-reserved	non-reserved	non-reserved	
COLUMN_NAME		non-reserved	non-reserved	non-reserved
COMMAND_FUNCTION		non-reserved	non-reserved	non-reserved
COMMAND_FUNCTION_CODE		non-reserved	non-reserved	
COMMENT	non-reserved			
COMMENTS	non-reserved			
COMMIT	non-reserved	reserved	reserved	reserved
COMMITTED	non-reserved	non-reserved	non-reserved	non-reserved
COMPRESSION	non-reserved			

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
CONCURRENTLY	reserved (can be function or type)			
CONDITION		reserved	reserved	
CONDITIONAL	non-reserved	non-reserved	non-reserved	
CONDITION_NUMBER		non-reserved	non-reserved	non-reserved
CONFIGURATION	non-reserved			
CONFLICT	non-reserved			
CONNECT		reserved	reserved	reserved
CONNECTION	non-reserved	non-reserved	non-reserved	reserved
CONNECTION_NAME		non-reserved	non-reserved	non-reserved
CONSTANT	non-reserved, requires AS			
CONSTRAINT	reserved	reserved	reserved	reserved
CONSTRAINTS	non-reserved	non-reserved	non-reserved	reserved
CONSTRAINT_CATALOG		non-reserved	non-reserved	non-reserved
CONSTRAINT_NAME		non-reserved	non-reserved	non-reserved
CONSTRAINT_SCHEMA		non-reserved	non-reserved	non-reserved
CONSTRUCTOR		non-reserved	non-reserved	
CONTAINS		reserved	reserved	
CONTENT	non-reserved	non-reserved	non-reserved	
CONTINUE	non-reserved	non-reserved	non-reserved	reserved
CONTROL		non-reserved	non-reserved	
CONVERSION	non-reserved			
CONVERT		reserved	reserved	reserved
COPARTITION		non-reserved		
COPY	non-reserved	reserved	reserved	
CORR		reserved	reserved	
CORRESPONDING		reserved	reserved	reserved
COS		reserved	reserved	
COSH		reserved	reserved	
COST	non-reserved			
COUNT		reserved	reserved	reserved
COVAR_POP		reserved	reserved	
COVAR_SAMP		reserved	reserved	
CREATE	reserved, requires AS	reserved	reserved	reserved
CROSS	reserved (can be function or type)	reserved	reserved	reserved
CSV	non-reserved			
CUBE	non-reserved	reserved	reserved	
CUME_DIST		reserved	reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
CURRENT	non-reserved	reserved	reserved	reserved
CURRENT_CATALOG	reserved	reserved	reserved	
CURRENT_DATE	reserved	reserved	reserved	reserved
CURRENT_DEFAULT_TRANSFORM_GROUP		reserved	reserved	
CURRENT_PATH		reserved	reserved	
CURRENT_ROLE	reserved	reserved	reserved	
CURRENT_ROW		reserved	reserved	
CURRENT_SCHEMA	reserved (can be function or type)	reserved	reserved	
CURRENT_TIME	reserved	reserved	reserved	reserved
CURRENT_TIMESTAMP	reserved	reserved	reserved	reserved
CURRENT_TRANSFORM_GROUP_FOR_TYPE		reserved	reserved	
CURRENT_USER	reserved	reserved	reserved	reserved
CURSOR	non-reserved	reserved	reserved	reserved
CURSOR_NAME		non-reserved	non-reserved	non-reserved
CYCLE	non-reserved	reserved	reserved	
DATA	non-reserved	non-reserved	non-reserved	non-reserved
DATABASE	non-reserved			
DATALINK		reserved	reserved	
DATE		reserved	reserved	reserved
DATETIME_INTERVAL_CODE		non-reserved	non-reserved	non-reserved
DATETIME_INTERVAL_PRECISION		non-reserved	non-reserved	non-reserved
DAY	non-reserved, requires AS	reserved	reserved	reserved
DB		non-reserved	non-reserved	
DEALLOCATE	non-reserved	reserved	reserved	reserved
DEC	non-reserved (cannot be function or type)	reserved	reserved	reserved
DECFLOAT		reserved	reserved	
DECIMAL	non-reserved (cannot be function or type)	reserved	reserved	reserved
DECLARE	non-reserved	reserved	reserved	reserved
DEFAULT	reserved	reserved	reserved	reserved
DEFAULTS	non-reserved	non-reserved	non-reserved	
DEFERRABLE	reserved	non-reserved	non-reserved	reserved
DEFERRED	non-reserved	non-reserved	non-reserved	reserved
DEFINE		reserved	reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
DEFINED		non-reserved	non-reserved	
DEFINER	non-reserved	non-reserved	non-reserved	
DEGREE		non-reserved	non-reserved	
DELETE	non-reserved	reserved	reserved	reserved
DELIMITER	non-reserved			
DELIMITERS	non-reserved			
DENSE_RANK		reserved	reserved	
DEPENDS	non-reserved			
DEPTH	non-reserved	non-reserved	non-reserved	
DEREF		reserved	reserved	
DERIVED		non-reserved	non-reserved	
DESC	reserved	non-reserved	non-reserved	reserved
DESCRIBE		reserved	reserved	reserved
DESCRIPTOR		non-reserved	non-reserved	reserved
DETACH	non-reserved			
DETERMINISTIC		reserved	reserved	
DIAGNOSTICS		non-reserved	non-reserved	reserved
DICTIONARY	non-reserved			
DISABLE	non-reserved			
DISCARD	non-reserved			
DISCONNECT		reserved	reserved	reserved
DISPATCH		non-reserved	non-reserved	
DISTINCT	reserved	reserved	reserved	reserved
DLNEWCOPY		reserved	reserved	
DLPREVIOUSCOPY		reserved	reserved	
DLURLCOMPLETE		reserved	reserved	
DLURLCOMPLETEONLY		reserved	reserved	
DLURLCOMPLETEWRITE		reserved	reserved	
DLURLPATH		reserved	reserved	
DLURLPATHONLY		reserved	reserved	
DLURLPATHWRITE		reserved	reserved	
DLURLSCHEME		reserved	reserved	
DLURLSERVER		reserved	reserved	
DLVALUE		reserved	reserved	
DO	reserved			
DOCUMENT	non-reserved	non-reserved	non-reserved	
DOMAIN	non-reserved	non-reserved	non-reserved	reserved
DOUBLE	non-reserved	reserved	reserved	reserved
DROP	non-reserved	reserved	reserved	reserved
DYNAMIC		reserved	reserved	

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
DYNAMIC_FUNCTION		non-reserved	non-reserved	non-reserved
DYNAMIC_FUNCTION_CODE		non-reserved	non-reserved	
EACH	non-reserved	reserved	reserved	
ELEMENT		reserved	reserved	
ELSE	reserved	reserved	reserved	reserved
EMPTY	non-reserved	reserved	reserved	
ENABLE	non-reserved			
ENCODING	non-reserved	non-reserved	non-reserved	
ENCRYPTED	non-reserved			
END	reserved	reserved	reserved	reserved
END-EXEC		reserved	reserved	reserved
END_FRAME		reserved	reserved	
END_PARTITION		reserved	reserved	
ENFORCED		non-reserved	non-reserved	
ENUM	non-reserved			
EQUALS		reserved	reserved	
ERROR	non-reserved	non-reserved	non-reserved	
ESCAPE	non-reserved	reserved	reserved	reserved
EVENT	non-reserved			
EVERY		reserved	reserved	
EXCEPT	reserved, requires AS	reserved	reserved	reserved
EXCEPTION				reserved
EXCLUDE	non-reserved	non-reserved	non-reserved	
EXCLUDING	non-reserved	non-reserved	non-reserved	
EXCLUSIVE	non-reserved			
EXEC		reserved	reserved	reserved
EXECUTE	non-reserved	reserved	reserved	reserved
EXISTS	non-reserved (can- not be function or type)	reserved	reserved	reserved
EXP		reserved	reserved	
EXPLAIN	non-reserved			
EXPRESSION	non-reserved	non-reserved	non-reserved	
EXTENSION	non-reserved			
EXTERNAL	non-reserved	reserved	reserved	reserved
EXTRACT	non-reserved (can- not be function or type)	reserved	reserved	reserved
FAILED_AUTH_KEEP_TIME	non-reserved			
FAILED_LOGIN_ATTEMPTS	non-reserved			

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
FALSE	reserved	reserved	reserved	reserved
FAMILY	non-reserved			
FETCH	reserved, requires AS	reserved	reserved	reserved
FILE		non-reserved	non-reserved	
FILTER	non-reserved, requires AS	reserved	reserved	
FINAL		non-reserved	non-reserved	
FINALIZE	non-reserved			
FINISH		non-reserved	non-reserved	
FIRST	non-reserved	non-reserved	non-reserved	reserved
FIRST_VALUE		reserved	reserved	
FLAG		non-reserved	non-reserved	
FLOAT	non-reserved (cannot be function or type)	reserved	reserved	reserved
FLOOR		reserved	reserved	
FOLLOWING	non-reserved	non-reserved	non-reserved	
FOR	reserved, requires AS	reserved	reserved	reserved
FORCE	non-reserved			
FOREIGN	reserved	reserved	reserved	reserved
FORMAT	non-reserved	non-reserved	non-reserved	
FORTRAN		non-reserved	non-reserved	non-reserved
FORWARD	non-reserved			
FOUND		non-reserved	non-reserved	reserved
FRAME_ROW		reserved	reserved	
FREE		reserved	reserved	
FREEZE	reserved (can be function or type)			
FROM	reserved, requires AS	reserved	reserved	reserved
FS		non-reserved	non-reserved	
FULFILL		non-reserved	non-reserved	
FULL	reserved (can be function or type)	reserved	reserved	reserved
FUNCTION	non-reserved	reserved	reserved	
FUNCTIONS	non-reserved			
FUSION		reserved	reserved	
G		non-reserved	non-reserved	
GENERAL		non-reserved	non-reserved	
GENERATED	non-reserved	non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
GET		reserved	reserved	reserved
GLOBAL	non-reserved	reserved	reserved	reserved
GO		non-reserved	non-reserved	reserved
GOTO		non-reserved	non-reserved	reserved
GRANT	reserved, requires AS	reserved	reserved	reserved
GRANTED	non-reserved	non-reserved	non-reserved	
GREATEST	non-reserved (cannot be function or type)	reserved		
GROUP	reserved, requires AS	reserved	reserved	reserved
GROUPING	non-reserved (cannot be function or type)	reserved	reserved	
GROUPS	non-reserved	reserved	reserved	
HANDLER	non-reserved			
HAVING	reserved, requires AS	reserved	reserved	reserved
HEADER	non-reserved			
HEX		non-reserved	non-reserved	
HIERARCHY		non-reserved	non-reserved	
HOLD	non-reserved	reserved	reserved	
HOURL	non-reserved, requires AS	reserved	reserved	reserved
ID		non-reserved	non-reserved	
IDENTITY	non-reserved	reserved	reserved	reserved
IF	non-reserved			
IGNORE		non-reserved	non-reserved	
ILIKE	reserved (can be function or type)			
IMMEDIATE	non-reserved	non-reserved	non-reserved	reserved
IMMEDIATELY		non-reserved	non-reserved	
IMMUTABLE	non-reserved			
IMPLEMENTATION		non-reserved	non-reserved	
IMPLICIT	non-reserved			
IMPORT	non-reserved	reserved	reserved	
IN	reserved	reserved	reserved	reserved
INCLUDE	non-reserved			
INCLUDING	non-reserved	non-reserved	non-reserved	
INCREMENT	non-reserved	non-reserved	non-reserved	
INDENT	non-reserved	non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
INDEX	non-reserved			
INDEXES	non-reserved			
INDICATOR		reserved	reserved	reserved
INHERIT	non-reserved			
INHERITS	non-reserved			
INITIAL		reserved	reserved	
INITIALLY	reserved	non-reserved	non-reserved	reserved
INLINE	non-reserved			
INNER	reserved (can be function or type)	reserved	reserved	reserved
INOUT	non-reserved (cannot be function or type)	reserved	reserved	
INPUT	non-reserved	non-reserved	non-reserved	reserved
INSENSITIVE	non-reserved	reserved	reserved	reserved
INSERT	non-reserved	reserved	reserved	reserved
INSTANCE		non-reserved	non-reserved	
INSTANTIABLE		non-reserved	non-reserved	
INSTEAD	non-reserved	non-reserved	non-reserved	
INT	non-reserved (cannot be function or type)	reserved	reserved	reserved
INTEGER	non-reserved (cannot be function or type)	reserved	reserved	reserved
INTEGRITY		non-reserved	non-reserved	
INTERSECT	reserved, requires AS	reserved	reserved	reserved
INTERSECTION		reserved	reserved	
INTERVAL	non-reserved (cannot be function or type)	reserved	reserved	reserved
INTO	reserved, requires AS	reserved	reserved	reserved
INVOKER	non-reserved	non-reserved	non-reserved	
IS	reserved (can be function or type)	reserved	reserved	reserved
ISNULL	reserved (can be function or type), requires AS			
ISOLATION	non-reserved	non-reserved	non-reserved	reserved
JOIN	reserved (can be function or type)	reserved	reserved	reserved

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
JSON	non-reserved (cannot be function or type)	reserved		
JSON_ARRAY	non-reserved (cannot be function or type)	reserved	reserved	
JSON_ARRAYAGG	non-reserved (cannot be function or type)	reserved	reserved	
JSON_EXISTS	non-reserved (cannot be function or type)	reserved	reserved	
JSON_OBJECT	non-reserved (cannot be function or type)	reserved	reserved	
JSON_OBJECTAGG	non-reserved (cannot be function or type)	reserved	reserved	
JSON_QUERY	non-reserved (cannot be function or type)	reserved	reserved	
JSON_SCALAR	non-reserved (cannot be function or type)	reserved		
JSON_SERIALIZE	non-reserved (cannot be function or type)	reserved		
JSON_TABLE	non-reserved (cannot be function or type)	reserved	reserved	
JSON_TABLE_PRIMITIVE		reserved	reserved	
JSON_VALUE	non-reserved (cannot be function or type)	reserved	reserved	
K		non-reserved	non-reserved	
KEEP	non-reserved	non-reserved	non-reserved	
KEY	non-reserved	non-reserved	non-reserved	reserved
KEYS	non-reserved	non-reserved	non-reserved	
KEY_MEMBER		non-reserved	non-reserved	
KEY_TYPE		non-reserved	non-reserved	
LABEL	non-reserved			
LAG		reserved	reserved	
LANGUAGE	non-reserved	reserved	reserved	reserved
LARGE	non-reserved	reserved	reserved	
LAST	non-reserved	non-reserved	non-reserved	reserved
LAST_VALUE		reserved	reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
LATERAL	reserved	reserved	reserved	
LEAD		reserved	reserved	
LEADING	reserved	reserved	reserved	reserved
LEAKPROOF	non-reserved			
LEAST	non-reserved (cannot be function or type)	reserved		
LEFT	reserved (can be function or type)	reserved	reserved	reserved
LENGTH		non-reserved	non-reserved	non-reserved
LEVEL	non-reserved	non-reserved	non-reserved	reserved
LIBRARY		non-reserved	non-reserved	
LIKE	reserved (can be function or type)	reserved	reserved	reserved
LIKE_REGEX		reserved	reserved	
LIMIT	reserved, requires AS	non-reserved	non-reserved	
LINK		non-reserved	non-reserved	
LISTAGG		reserved	reserved	
LISTEN	non-reserved			
LN		reserved	reserved	
LOAD	non-reserved			
LOCAL	non-reserved	reserved	reserved	reserved
LOCALTIME	reserved	reserved	reserved	
LOCALTIMESTAMP	reserved	reserved	reserved	
LOCATION	non-reserved	non-reserved	non-reserved	
LOCATOR		non-reserved	non-reserved	
LOCK	non-reserved			
LOCKED	non-reserved			
LOG		reserved	reserved	
LOG10		reserved	reserved	
LOGGED	non-reserved			
LOWER		reserved	reserved	reserved
LPAD		reserved		
LTRIM		reserved		
M		non-reserved	non-reserved	
MAP		non-reserved	non-reserved	
MAPPING	non-reserved	non-reserved	non-reserved	
MATCH	non-reserved	reserved	reserved	reserved
MATCHED	non-reserved	non-reserved	non-reserved	
MATCHES		reserved	reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
MATCH_NUMBER		reserved	reserved	
MATCH_RECOGNIZE		reserved	reserved	
MATERIALIZED	non-reserved			
MAX		reserved	reserved	reserved
MAXVALUE	non-reserved	non-reserved	non-reserved	
MEASURES		non-reserved	non-reserved	
MEMBER		reserved	reserved	
MERGE	non-reserved	reserved	reserved	
MERGE_ACTION	non-reserved (cannot be function or type)			
MESSAGE_LENGTH		non-reserved	non-reserved	non-reserved
MESSAGE_OCTET_LENGTH		non-reserved	non-reserved	non-reserved
MESSAGE_TEXT		non-reserved	non-reserved	non-reserved
METHOD	non-reserved	reserved	reserved	
MIN		reserved	reserved	reserved
MINUTE	non-reserved, requires AS	reserved	reserved	reserved
MINVALUE	non-reserved	non-reserved	non-reserved	
MOD		reserved	reserved	
MODE	non-reserved			
MODIFIES		reserved	reserved	
MODULE		reserved	reserved	reserved
MONTH	non-reserved, requires AS	reserved	reserved	reserved
MORE		non-reserved	non-reserved	non-reserved
MOVE	non-reserved			
MULTISET		reserved	reserved	
MUMPS		non-reserved	non-reserved	non-reserved
NAME	non-reserved	non-reserved	non-reserved	non-reserved
NAMES	non-reserved	non-reserved	non-reserved	reserved
NAMESPACE		non-reserved	non-reserved	
NATIONAL	non-reserved (cannot be function or type)	reserved	reserved	reserved
NATURAL	reserved (can be function or type)	reserved	reserved	reserved
NCHAR	non-reserved (cannot be function or type)	reserved	reserved	reserved
NCLOB		reserved	reserved	
NESTED	non-reserved	non-reserved	non-reserved	
NESTING		non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
NEW	non-reserved	reserved	reserved	
NEXT	non-reserved	non-reserved	non-reserved	reserved
NFC	non-reserved	non-reserved	non-reserved	
NFD	non-reserved	non-reserved	non-reserved	
NFKC	non-reserved	non-reserved	non-reserved	
NFKD	non-reserved	non-reserved	non-reserved	
NIL		non-reserved	non-reserved	
NO	non-reserved	reserved	reserved	reserved
NONE	non-reserved (cannot be function or type)	reserved	reserved	
NORMALIZE	non-reserved (cannot be function or type)	reserved	reserved	
NORMALIZED	non-reserved	non-reserved	non-reserved	
NOT	reserved	reserved	reserved	reserved
NOTHING	non-reserved			
NOTIFY	non-reserved			
NOTNULL	reserved (can be function or type), requires AS			
NOWAIT	non-reserved			
NTH_VALUE		reserved	reserved	
NTILE		reserved	reserved	
NULL	reserved	reserved	reserved	reserved
NULLABLE		non-reserved	non-reserved	non-reserved
NULLIF	non-reserved (cannot be function or type)	reserved	reserved	reserved
NULLS	non-reserved	non-reserved	non-reserved	
NULL_ORDERING		non-reserved	non-reserved	
NUMBER		non-reserved	non-reserved	non-reserved
NUMERIC	non-reserved (cannot be function or type)	reserved	reserved	reserved
OBJECT	non-reserved	non-reserved	non-reserved	
OCCURRENCE		non-reserved	non-reserved	
OCCURRENCES_REGEX		reserved	reserved	
OCTETS		non-reserved	non-reserved	
OCTET_LENGTH		reserved	reserved	reserved
OF	non-reserved	reserved	reserved	reserved
OFF	non-reserved	non-reserved	non-reserved	
OFFICER	non-reserved			

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
OFFSET	reserved, requires AS	reserved	reserved	
OIDS	non-reserved			
OLD	non-reserved	reserved	reserved	
OMIT	non-reserved	reserved	reserved	
ON	reserved, requires AS	reserved	reserved	reserved
ONE		reserved	reserved	
ONLY	reserved	reserved	reserved	reserved
OPEN		reserved	reserved	reserved
OPERATOR	non-reserved			
OPTION	non-reserved	non-reserved	non-reserved	reserved
OPTIONS	non-reserved	non-reserved	non-reserved	
OR	reserved	reserved	reserved	reserved
ORDER	reserved, requires AS	reserved	reserved	reserved
ORDERING		non-reserved	non-reserved	
ORDINALITY	non-reserved	non-reserved	non-reserved	
OTHERS	non-reserved	non-reserved	non-reserved	
OUT	non-reserved (cannot be function or type)	reserved	reserved	
OUTER	reserved (can be function or type)	reserved	reserved	reserved
OUTPUT		non-reserved	non-reserved	reserved
OVER	non-reserved, requires AS	reserved	reserved	
OVERFLOW		non-reserved	non-reserved	
OVERLAPS	reserved (can be function or type), requires AS	reserved	reserved	reserved
OVERLAY	non-reserved (cannot be function or type)	reserved	reserved	
OVERRIDING	non-reserved	non-reserved	non-reserved	
OWNED	non-reserved			
OWNER	non-reserved			
P		non-reserved	non-reserved	
PACKAGE	non-reserved			
PAD		non-reserved	non-reserved	reserved
PARALLEL	non-reserved			
PARAMETER	non-reserved	reserved	reserved	
PARAMETER_MODE		non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
PARAMETER_NAME		non-reserved	non-reserved	
PARAMETER_ORDINAL_POSITION		non-reserved	non-reserved	
PARAMETER_SPECIFIC_CATALOG		non-reserved	non-reserved	
PARAMETER_SPECIFIC_NAME		non-reserved	non-reserved	
PARAMETER_SPECIFIC_SCHEMA		non-reserved	non-reserved	
PARSER	non-reserved			
PARTIAL	non-reserved	non-reserved	non-reserved	reserved
PARTIAL_AGGREGATE	reserved			
PARTITION	non-reserved	reserved	reserved	
PARTITIONS	non-reserved			
PASCAL		non-reserved	non-reserved	non-reserved
PASS		non-reserved	non-reserved	
PASSING	non-reserved	non-reserved	non-reserved	
PASSTHROUGH		non-reserved	non-reserved	
PASSWORD	non-reserved			
PASSWORD_GRACE_TIME	non-reserved			
PASSWORD_LIFE_TIME	non-reserved			
PASSWORD_MIN_LEN	non-reserved			
PASSWORD_MIN_UNIQUE_CHARS	non-reserved			
PASSWORD_REQUIRE_COMPLEX	non-reserved			
PASSWORD_REUSE_MAX	non-reserved			
PASSWORD_REUSE_TIME	non-reserved			
PAST		non-reserved	non-reserved	
PATH	non-reserved	non-reserved	non-reserved	
PATTERN		reserved	reserved	
PER		reserved	reserved	
PERCENT		reserved	reserved	
PERCENTILE_CONT		reserved	reserved	
PERCENTILE_DISC		reserved	reserved	
PERCENT_RANK		reserved	reserved	
PERIOD		reserved	reserved	
PERMISSION		non-reserved	non-reserved	
PERMUTE		non-reserved	non-reserved	
PIPE		non-reserved	non-reserved	
PLACING	reserved	non-reserved	non-reserved	
PLAN	non-reserved	non-reserved	non-reserved	
PLANS	non-reserved			
PLI		non-reserved	non-reserved	non-reserved
POLICY	non-reserved			

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
PORTION		reserved	reserved	
POSITION	non-reserved (cannot be function or type)	reserved	reserved	reserved
POSITION_REGEX		reserved	reserved	
POWER		reserved	reserved	
PRECEDES		reserved	reserved	
PRECEDING	non-reserved	non-reserved	non-reserved	
PRECISION	non-reserved (cannot be function or type), requires AS	reserved	reserved	reserved
PREPARE	non-reserved	reserved	reserved	reserved
PREPARED	non-reserved			
PRESERVE	non-reserved	non-reserved	non-reserved	reserved
PREV		non-reserved	non-reserved	
PRIMARY	reserved	reserved	reserved	reserved
PRIOR	non-reserved	non-reserved	non-reserved	reserved
PRIVATE		non-reserved	non-reserved	
PRIVILEGES	non-reserved	non-reserved	non-reserved	reserved
PROCEDURAL	non-reserved			
PROCEDURE	non-reserved	reserved	reserved	reserved
PROCEDURES	non-reserved			
PROFILE	non-reserved			
PROGRAM	non-reserved			
PRUNE		non-reserved	non-reserved	
PTF		reserved	reserved	
PUBLIC		non-reserved	non-reserved	reserved
PUBLICATION	non-reserved			
QUOTE	non-reserved			
QUOTES	non-reserved	non-reserved	non-reserved	
RANGE	non-reserved	reserved	reserved	
RANK		reserved	reserved	
READ	non-reserved	non-reserved	non-reserved	reserved
READS		reserved	reserved	
REAL	non-reserved (cannot be function or type)	reserved	reserved	reserved
REASSIGN	non-reserved			
RECHECK	non-reserved			
RECOVERY		non-reserved	non-reserved	
RECURSIVE	non-reserved	reserved	reserved	
REF	non-reserved	reserved	reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
REFERENCES	reserved	reserved	reserved	reserved
REFERENCING	non-reserved	reserved	reserved	
REFRESH	non-reserved			
REGR_AVGX		reserved	reserved	
REGR_AVGY		reserved	reserved	
REGR_COUNT		reserved	reserved	
REGR_INTERCEPT		reserved	reserved	
REGR_R2		reserved	reserved	
REGR_SLOPE		reserved	reserved	
REGR_SXX		reserved	reserved	
REGR_SXY		reserved	reserved	
REGR_SYY		reserved	reserved	
REINDEX	non-reserved			
RELATIVE	non-reserved	non-reserved	non-reserved	reserved
RELEASE	non-reserved	reserved	reserved	
RENAME	non-reserved			
REPEATABLE	non-reserved	non-reserved	non-reserved	non-reserved
REPLACE	non-reserved			
REPLICA	non-reserved			
REQUIRING		non-reserved	non-reserved	
RESET	non-reserved			
RESPECT		non-reserved	non-reserved	
RESTART	non-reserved	non-reserved	non-reserved	
RESTORE		non-reserved	non-reserved	
RESTRICT	non-reserved	non-reserved	non-reserved	reserved
RESULT		reserved	reserved	
RETURN	non-reserved	reserved	reserved	
RETURNED_CARDINALITY		non-reserved	non-reserved	
RETURNED_LENGTH		non-reserved	non-reserved	non-reserved
RETURNED_OCTET_LENGTH		non-reserved	non-reserved	non-reserved
RETURNED_SQLSTATE		non-reserved	non-reserved	non-reserved
RETURNING	reserved, requires AS	non-reserved	non-reserved	
RETURNS	non-reserved	reserved	reserved	
REVOKE	non-reserved	reserved	reserved	reserved
RIGHT	reserved (can be function or type)	reserved	reserved	reserved
ROLE	non-reserved	non-reserved	non-reserved	
ROLLBACK	non-reserved	reserved	reserved	reserved
ROLLUP	non-reserved	reserved	reserved	

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
ROUTINE	non-reserved	non-reserved	non-reserved	
ROUTINES	non-reserved			
ROUTINE_CATALOG		non-reserved	non-reserved	
ROUTINE_NAME		non-reserved	non-reserved	
ROUTINE_SCHEMA		non-reserved	non-reserved	
ROW	non-reserved (cannot be function or type)	reserved	reserved	
ROWS	non-reserved	reserved	reserved	reserved
ROW_COUNT		non-reserved	non-reserved	non-reserved
ROW_NUMBER		reserved	reserved	
RPAD		reserved		
RTRIM		reserved		
RULE	non-reserved			
RUNNING		reserved	reserved	
SAVEPOINT	non-reserved	reserved	reserved	
SCALAR	non-reserved	non-reserved	non-reserved	
SCALE		non-reserved	non-reserved	non-reserved
SCHEMA	non-reserved	non-reserved	non-reserved	reserved
SCHEMAS	non-reserved			
SCHEMA_NAME		non-reserved	non-reserved	non-reserved
SCOPE		reserved	reserved	
SCOPE_CATALOG		non-reserved	non-reserved	
SCOPE_NAME		non-reserved	non-reserved	
SCOPE_SCHEMA		non-reserved	non-reserved	
SCROLL	non-reserved	reserved	reserved	reserved
SEARCH	non-reserved	reserved	reserved	
SECOND	non-reserved, requires AS	reserved	reserved	reserved
SECTION		non-reserved	non-reserved	reserved
SECURITY	non-reserved	non-reserved	non-reserved	
SEEK		reserved	reserved	
SELECT	reserved	reserved	reserved	reserved
SELECTIVE		non-reserved	non-reserved	
SELF		non-reserved	non-reserved	
SEMANTICS		non-reserved	non-reserved	
SENSITIVE		reserved	reserved	
SEQUENCE	non-reserved	non-reserved	non-reserved	
SEQUENCES	non-reserved			
SERIALIZABLE	non-reserved	non-reserved	non-reserved	non-reserved
SERVER	non-reserved	non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
SERVER_NAME		non-reserved	non-reserved	non-reserved
SESSION	non-reserved	non-reserved	non-reserved	reserved
SESSION_USER	reserved	reserved	reserved	reserved
SET	non-reserved	reserved	reserved	reserved
SETOF	non-reserved (cannot be function or type)			
SETS	non-reserved	non-reserved	non-reserved	
SHARE	non-reserved			
SHOW	non-reserved	reserved	reserved	
SIMILAR	reserved (can be function or type)	reserved	reserved	
SIMPLE	non-reserved	non-reserved	non-reserved	
SIN		reserved	reserved	
SINH		reserved	reserved	
SIZE		non-reserved	non-reserved	reserved
SKIP	non-reserved	reserved	reserved	
SMALLINT	non-reserved (cannot be function or type)	reserved	reserved	reserved
SNAPSHOT	non-reserved			
SOME	reserved	reserved	reserved	reserved
SORT_DIRECTION		non-reserved	non-reserved	
SOURCE	non-reserved	non-reserved	non-reserved	
SPACE		non-reserved	non-reserved	reserved
SPECIFIC		reserved	reserved	
SPECIFICTYPE		reserved	reserved	
SPECIFIC_NAME		non-reserved	non-reserved	
SPLIT	non-reserved			
SQL	non-reserved	reserved	reserved	reserved
SQLCODE				reserved
SQLERROR				reserved
SQLEXCEPTION		reserved	reserved	
SQLSTATE		reserved	reserved	reserved
SQLWARNING		reserved	reserved	
SQRT		reserved	reserved	
STABLE	non-reserved			
STANDALONE	non-reserved	non-reserved	non-reserved	
START	non-reserved	reserved	reserved	
STATE		non-reserved	non-reserved	
STATEMENT	non-reserved	non-reserved	non-reserved	

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
STATIC		reserved	reserved	
STATISTICS	non-reserved			
STDDEV_POP		reserved	reserved	
STDDEV_SAMP		reserved	reserved	
STDIN	non-reserved			
STDOUT	non-reserved			
STORAGE	non-reserved			
STORED	non-reserved			
STRICT	non-reserved			
STRING	non-reserved	non-reserved	non-reserved	
STRIP	non-reserved	non-reserved	non-reserved	
STRUCTURE		non-reserved	non-reserved	
STYLE		non-reserved	non-reserved	
SUBCLASS_ORIGIN		non-reserved	non-reserved	non-reserved
SUBMULTISET		reserved	reserved	
SUBSCRIPTION	non-reserved			
SUBSET		reserved	reserved	
SUBSTRING	non-reserved (cannot be function or type)	reserved	reserved	reserved
SUBSTRING_REGEX		reserved	reserved	
SUCCEEDS		reserved	reserved	
SUM		reserved	reserved	reserved
SUPPORT	non-reserved			
SYMMETRIC	reserved	reserved	reserved	
SYSID	non-reserved			
SYSTEM	non-reserved	reserved	reserved	
SYSTEM_TIME		reserved	reserved	
SYSTEM_USER	reserved	reserved	reserved	reserved
T		non-reserved	non-reserved	
TABLE	reserved	reserved	reserved	reserved
TABLES	non-reserved			
TABLESAMPLE	reserved (can be function or type)	reserved	reserved	
TABLESPACE	non-reserved			
TABLE_NAME		non-reserved	non-reserved	non-reserved
TAN		reserved	reserved	
TANH		reserved	reserved	
TARGET	non-reserved			
TEMP	non-reserved			
TEMPLATE	non-reserved			

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
TEMPORARY	non-reserved	non-reserved	non-reserved	reserved
TEXT	non-reserved			
THEN	reserved	reserved	reserved	reserved
THROUGH		non-reserved	non-reserved	
TIES	non-reserved	non-reserved	non-reserved	
TIME	non-reserved (cannot be function or type)	reserved	reserved	reserved
TIMESTAMP	non-reserved (cannot be function or type)	reserved	reserved	reserved
TIMEZONE_HOUR		reserved	reserved	reserved
TIMEZONE_MINUTE		reserved	reserved	reserved
TO	reserved, requires AS	reserved	reserved	reserved
TOKEN		non-reserved	non-reserved	
TOP_LEVEL_COUNT		non-reserved	non-reserved	
TRAILING	reserved	reserved	reserved	reserved
TRANSACTION	non-reserved	non-reserved	non-reserved	reserved
TRANSACTIONS_COMMITTED		non-reserved	non-reserved	
TRANSACTIONS_ROLLED_BACK		non-reserved	non-reserved	
TRANSACTION_ACTIVE		non-reserved	non-reserved	
TRANSFORM	non-reserved	non-reserved	non-reserved	
TRANSFORMS		non-reserved	non-reserved	
TRANSLATE		reserved	reserved	reserved
TRANSLATE_REGEX		reserved	reserved	
TRANSLATION		reserved	reserved	reserved
TREAT	non-reserved (cannot be function or type)	reserved	reserved	
TRIGGER	non-reserved	reserved	reserved	
TRIGGER_CATALOG		non-reserved	non-reserved	
TRIGGER_NAME		non-reserved	non-reserved	
TRIGGER_SCHEMA		non-reserved	non-reserved	
TRIM	non-reserved (cannot be function or type)	reserved	reserved	reserved
TRIM_ARRAY		reserved	reserved	
TRUE	reserved	reserved	reserved	reserved
TRUNCATE	non-reserved	reserved	reserved	
TRUSTED	non-reserved			
TYPE	non-reserved	non-reserved	non-reserved	non-reserved
TYPES	non-reserved			

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
UESCAPE	non-reserved	reserved	reserved	
UNBOUNDED	non-reserved	non-reserved	non-reserved	
UNCOMMITTED	non-reserved	non-reserved	non-reserved	non-reserved
UNCONDITIONAL	non-reserved	non-reserved	non-reserved	
UNDER		non-reserved	non-reserved	
UNENCRYPTED	non-reserved			
UNION	reserved, requires AS	reserved	reserved	reserved
UNIQUE	reserved	reserved	reserved	reserved
UNKNOWN	non-reserved	reserved	reserved	reserved
UNLIMITED	non-reserved			
UNLINK		non-reserved	non-reserved	
UNLISTEN	non-reserved			
UNLOCK	non-reserved			
UNLOGGED	non-reserved			
UNMATCHED		non-reserved	non-reserved	
UNNAMED		non-reserved	non-reserved	non-reserved
UNNEST		reserved	reserved	
UNTIL	non-reserved			
UNTYPED		non-reserved	non-reserved	
UPDATE	non-reserved	reserved	reserved	reserved
UPPER		reserved	reserved	reserved
URI		non-reserved	non-reserved	
USAGE		non-reserved	non-reserved	reserved
USER	reserved	reserved	reserved	reserved
USER_DEFINED_TYPE_CATALOG		non-reserved	non-reserved	
USER_DEFINED_TYPE_CODE		non-reserved	non-reserved	
USER_DEFINED_TYPE_NAME		non-reserved	non-reserved	
USER_DEFINED_TYPE_SCHEMA		non-reserved	non-reserved	
USER_INACTIVE_TIME	non-reserved			
USING	reserved	reserved	reserved	reserved
UTF16		non-reserved	non-reserved	
UTF32		non-reserved	non-reserved	
UTF8		non-reserved	non-reserved	
VACUUM	non-reserved			
VALID	non-reserved	non-reserved	non-reserved	
VALIDATE	non-reserved			
VALIDATOR	non-reserved			
VALUE	non-reserved	reserved	reserved	reserved

SQL Key Words

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
VALUES	non-reserved (cannot be function or type)	reserved	reserved	reserved
VALUE_OF		reserved	reserved	
VARBINARY		reserved	reserved	
VARCHAR	non-reserved (cannot be function or type)	reserved	reserved	reserved
VARIADIC	reserved			
VARYING	non-reserved, requires AS	reserved	reserved	reserved
VAR_POP		reserved	reserved	
VAR_SAMP		reserved	reserved	
VERBOSE	reserved (can be function or type)			
VERSION	non-reserved	non-reserved	non-reserved	
VERSIONING		reserved	reserved	
VIEW	non-reserved	non-reserved	non-reserved	reserved
VIEWS	non-reserved			
VOLATILE	non-reserved			
WHEN	reserved	reserved	reserved	reserved
WHENEVER		reserved	reserved	reserved
WHERE	reserved, requires AS	reserved	reserved	reserved
WHITESPACE	non-reserved	non-reserved	non-reserved	
WIDTH_BUCKET		reserved	reserved	
WINDOW	reserved, requires AS	reserved	reserved	
WITH	reserved, requires AS	reserved	reserved	reserved
WITHIN	non-reserved, requires AS	reserved	reserved	
WITHOUT	non-reserved, requires AS	reserved	reserved	
WORK	non-reserved	non-reserved	non-reserved	reserved
WRAPPER	non-reserved	non-reserved	non-reserved	
WRITE	non-reserved	non-reserved	non-reserved	reserved
XML	non-reserved	reserved	reserved	
XMLAGG		reserved	reserved	
XMLATTRIBUTES	non-reserved (cannot be function or type)	reserved	reserved	
XMLBINARY		reserved	reserved	
XMLCAST		reserved	reserved	

Key Word	PostgreSQL	SQL:2023	SQL:2016	SQL-92
XMLCOMMENT		reserved	reserved	
XMLCONCAT	non-reserved (cannot be function or type)	reserved	reserved	
XMLDECLARATION		non-reserved	non-reserved	
XMLDOCUMENT		reserved	reserved	
XMLELEMENT	non-reserved (cannot be function or type)	reserved	reserved	
XML EXISTS	non-reserved (cannot be function or type)	reserved	reserved	
XMLFOREST	non-reserved (cannot be function or type)	reserved	reserved	
XMLITERATE		reserved	reserved	
XMLNAMESPACES	non-reserved (cannot be function or type)	reserved	reserved	
XMLPARSE	non-reserved (cannot be function or type)	reserved	reserved	
XMLPI	non-reserved (cannot be function or type)	reserved	reserved	
XMLQUERY		reserved	reserved	
XMLROOT	non-reserved (cannot be function or type)			
XMLSCHEMA		non-reserved	non-reserved	
XMLSERIALIZE	non-reserved (cannot be function or type)	reserved	reserved	
XMLTABLE	non-reserved (cannot be function or type)	reserved	reserved	
XMLTEXT		reserved	reserved	
XMLVALIDATE		reserved	reserved	
YEAR	non-reserved, requires AS	reserved	reserved	reserved
YES	non-reserved	non-reserved	non-reserved	
ZONE	non-reserved	non-reserved	non-reserved	reserved

Appendix D. SQL Conformance

This section attempts to outline to what extent Postgres Pro conforms to the current SQL standard. The following information is not a full statement of conformance, but it presents the main topics in as much detail as is both reasonable and useful for users.

The formal name of the SQL standard is ISO/IEC 9075 “Database Language SQL”. A revised version of the standard is released from time to time; the most recent update appearing in 2023. The 2023 version is referred to as ISO/IEC 9075:2023, or simply as SQL:2023. The versions prior to that were SQL:2016, SQL:2011, SQL:2008, SQL:2006, SQL:2003, SQL:1999, and SQL-92. Each version replaces the previous one, so claims of conformance to earlier versions have no official merit. Postgres Pro development aims for conformance with the latest official version of the standard where such conformance does not contradict traditional features or common sense. Many of the features required by the SQL standard are supported, though sometimes with slightly differing syntax or function. Further moves towards conformance can be expected over time.

SQL-92 defined three feature sets for conformance: Entry, Intermediate, and Full. Most database management systems claiming SQL standard conformance were conforming at only the Entry level, since the entire set of features in the Intermediate and Full levels was either too voluminous or in conflict with legacy behaviors.

Starting with SQL:1999, the SQL standard defines a large set of individual features rather than the ineffectively broad three levels found in SQL-92. A large subset of these features represents the “Core” features, which every conforming SQL implementation must supply. The rest of the features are purely optional.

The standard is split into a number of parts, each also known by a shorthand name:

- ISO/IEC 9075-1 Framework (SQL/Framework)
- ISO/IEC 9075-2 Foundation (SQL/Foundation)
- ISO/IEC 9075-3 Call Level Interface (SQL/CLI)
- ISO/IEC 9075-4 Persistent Stored Modules (SQL/PSM)
- ISO/IEC 9075-9 Management of External Data (SQL/MED)
- ISO/IEC 9075-10 Object Language Bindings (SQL/OLB)
- ISO/IEC 9075-11 Information and Definition Schemas (SQL/Schemata)
- ISO/IEC 9075-13 Routines and Types using the Java Language (SQL/JRT)
- ISO/IEC 9075-14 XML-related specifications (SQL/XML)
- ISO/IEC 9075-15 Multi-dimensional arrays (SQL/MDA)
- ISO/IEC 9075-16 Property Graph Queries (SQL/PGQ)

Note that some part numbers are not (or no longer) used.

The Postgres Pro core covers parts 1, 2, 9, 11, and 14. Part 3 is covered by the ODBC driver, and part 13 is covered by the PL/Java plug-in, but exact conformance is currently not being verified for these components. There are currently no implementations of parts 4, 10, 15, and 16 for Postgres Pro.

Postgres Pro supports most of the major features of SQL:2023. Out of 177 mandatory features required for full Core conformance, Postgres Pro conforms to at least 170. In addition, there is a long list of

supported optional features. It might be worth noting that at the time of writing, no current version of any database management system claims full conformance to Core SQL:2023.

In the following two sections, we provide a list of those features that Postgres Pro supports, followed by a list of the features defined in SQL:2023 which are not yet supported in Postgres Pro. Both of these lists are approximate: There might be minor details that are nonconforming for a feature that is listed as supported, and large parts of an unsupported feature might in fact be implemented. The main body of the documentation always contains the most accurate information about what does and does not work.

Note

Feature codes containing a hyphen are subfeatures. Therefore, if a particular subfeature is not supported, the main feature is listed as unsupported even if some other subfeatures are supported.

D.1. Supported Features

Identifier	Core?	Description	Comment
B012		Embedded C	
B021		Direct SQL	
B128		Routine language SQL	
E011	Core	Numeric data types	
E011-01	Core	INTEGER and SMALLINT data types	
E011-02	Core	REAL, DOUBLE PRECISION, and FLOAT data types	
E011-03	Core	DECIMAL and NUMERIC data types	
E011-04	Core	Arithmetic operators	
E011-05	Core	Numeric comparison	
E011-06	Core	Implicit casting among the numeric data types	
E021	Core	Character data types	
E021-01	Core	CHARACTER data type	
E021-02	Core	CHARACTER VARYING data type	
E021-03	Core	Character literals	
E021-04	Core	CHARACTER_LENGTH function	trims trailing spaces from CHARACTER values before counting
E021-05	Core	OCTET_LENGTH function	
E021-06	Core	SUBSTRING function	
E021-07	Core	Character concatenation	
E021-08	Core	UPPER and LOWER functions	
E021-09	Core	TRIM function	
E021-10	Core	Implicit casting among the character string types	
E021-11	Core	POSITION function	
E021-12	Core	Character comparison	
E031	Core	Identifiers	
E031-01	Core	Delimited identifiers	
E031-02	Core	Lower case identifiers	
E031-03	Core	Trailing underscore	

Identifier	Core?	Description	Comment
E051	Core	Basic query specification	
E051-01	Core	SELECT DISTINCT	
E051-02	Core	GROUP BY clause	
E051-04	Core	GROUP BY can contain columns not in <select list>	
E051-05	Core	Select list items can be renamed	
E051-06	Core	HAVING clause	
E051-07	Core	Qualified * in select list	
E051-08	Core	Correlation names in the FROM clause	
E051-09	Core	Rename columns in the FROM clause	
E061	Core	Basic predicates and search conditions	
E061-01	Core	Comparison predicate	
E061-02	Core	BETWEEN predicate	
E061-03	Core	IN predicate with list of values	
E061-04	Core	LIKE predicate	
E061-05	Core	LIKE predicate ESCAPE clause	
E061-06	Core	NULL predicate	
E061-07	Core	Quantified comparison predicate	
E061-08	Core	EXISTS predicate	
E061-09	Core	Subqueries in comparison predicate	
E061-11	Core	Subqueries in IN predicate	
E061-12	Core	Subqueries in quantified comparison predicate	
E061-13	Core	Correlated subqueries	
E061-14	Core	Search condition	
E071	Core	Basic query expressions	
E071-01	Core	UNION DISTINCT table operator	
E071-02	Core	UNION ALL table operator	
E071-03	Core	EXCEPT DISTINCT table operator	
E071-05	Core	Columns combined via table operators need not have exactly the same data type	
E071-06	Core	Table operators in subqueries	
E081	Core	Basic Privileges	
E081-01	Core	SELECT privilege	
E081-02	Core	DELETE privilege	
E081-03	Core	INSERT privilege at the table level	
E081-04	Core	UPDATE privilege at the table level	
E081-05	Core	UPDATE privilege at the column level	
E081-06	Core	REFERENCES privilege at the table level	
E081-07	Core	REFERENCES privilege at the column level	
E081-08	Core	WITH GRANT OPTION	
E081-09	Core	USAGE privilege	
E081-10	Core	EXECUTE privilege	

Identifier	Core?	Description	Comment
E091	Core	Set functions	
E091-01	Core	AVG	
E091-02	Core	COUNT	
E091-03	Core	MAX	
E091-04	Core	MIN	
E091-05	Core	SUM	
E091-06	Core	ALL quantifier	
E091-07	Core	DISTINCT quantifier	
E101	Core	Basic data manipulation	
E101-01	Core	INSERT statement	
E101-03	Core	Searched UPDATE statement	
E101-04	Core	Searched DELETE statement	
E111	Core	Single row SELECT statement	
E121	Core	Basic cursor support	
E121-01	Core	DECLARE CURSOR	
E121-02	Core	ORDER BY columns need not be in select list	
E121-03	Core	Value expressions in ORDER BY clause	
E121-04	Core	OPEN statement	
E121-06	Core	Positioned UPDATE statement	
E121-07	Core	Positioned DELETE statement	
E121-08	Core	CLOSE statement	
E121-10	Core	FETCH statement implicit NEXT	
E121-17	Core	WITH HOLD cursors	
E131	Core	Null value support (nulls in lieu of values)	
E141	Core	Basic integrity constraints	
E141-01	Core	NOT NULL constraints	
E141-02	Core	UNIQUE constraints of NOT NULL columns	
E141-03	Core	PRIMARY KEY constraints	
E141-04	Core	Basic FOREIGN KEY constraint with the NO ACTION default for both referential delete action and referential update action	
E141-06	Core	CHECK constraints	
E141-07	Core	Column defaults	
E141-08	Core	NOT NULL inferred on PRIMARY KEY	
E141-10	Core	Names in a foreign key can be specified in any order	
E151	Core	Transaction support	
E151-01	Core	COMMIT statement	
E151-02	Core	ROLLBACK statement	
E152	Core	Basic SET TRANSACTION statement	
E152-01	Core	SET TRANSACTION statement: ISOLATION LEVEL SERIALIZABLE clause	

Identifier	Core?	Description	Comment
E152-02	Core	SET TRANSACTION statement: READ ONLY and READ WRITE clauses	
E153	Core	Updatable queries with subqueries	
E161	Core	SQL comments using leading double minus	
E171	Core	SQLSTATE support	
E182	Core	Host language binding	
F021	Core	Basic information schema	
F021-01	Core	COLUMNS view	
F021-02	Core	TABLES view	
F021-03	Core	VIEWS view	
F021-04	Core	TABLE_CONSTRAINTS view	
F021-05	Core	REFERENTIAL_CONSTRAINTS view	
F021-06	Core	CHECK_CONSTRAINTS view	
F031	Core	Basic schema manipulation	
F031-01	Core	CREATE TABLE statement to create persistent base tables	
F031-02	Core	CREATE VIEW statement	
F031-03	Core	GRANT statement	
F031-04	Core	ALTER TABLE statement: ADD COLUMN clause	
F031-13	Core	DROP TABLE statement: RESTRICT clause	
F031-16	Core	DROP VIEW statement: RESTRICT clause	
F031-19	Core	REVOKE statement: RESTRICT clause	
F032		CASCADE drop behavior	
F033		ALTER TABLE statement: DROP COLUMN clause	
F034		Extended REVOKE statement	
F035		REVOKE with CASCADE	
F036		REVOKE statement performed by non-owner	
F037		REVOKE statement: GRANT OPTION FOR clause	
F038		REVOKE of a WITH GRANT OPTION privilege	
F041	Core	Basic joined table	
F041-01	Core	Inner join (but not necessarily the INNER keyword)	
F041-02	Core	INNER keyword	
F041-03	Core	LEFT OUTER JOIN	
F041-04	Core	RIGHT OUTER JOIN	
F041-05	Core	Outer joins can be nested	
F041-07	Core	The inner table in a left or right outer join can also be used in an inner join	
F041-08	Core	All comparison operators are supported (rather than just =)	
F051	Core	Basic date and time	
F051-01	Core	DATE data type (including support of DATE literal)	

Identifier	Core?	Description	Comment
F051-02	Core	TIME data type (including support of TIME literal) with fractional seconds precision of at least 0	
F051-03	Core	TIMESTAMP data type (including support of TIMESTAMP literal) with fractional seconds precision of at least 0 and 6	
F051-04	Core	Comparison predicate on DATE, TIME, and TIMESTAMP data types	
F051-05	Core	Explicit CAST between datetime types and character string types	
F051-06	Core	CURRENT_DATE	
F051-07	Core	LOCALTIME	
F051-08	Core	LOCALTIMESTAMP	
F052		Intervals and datetime arithmetic	
F053		OVERLAPS predicate	
F081	Core	UNION and EXCEPT in views	
F111		Isolation levels other than SERIALIZABLE	
F112		Isolation level READ UNCOMMITTED	
F113		Isolation level READ COMMITTED	
F114		Isolation level REPEATABLE READ	
F131	Core	Grouped operations	
F131-01	Core	WHERE, GROUP BY, and HAVING clauses supported in queries with grouped views	
F131-02	Core	Multiple tables supported in queries with grouped views	
F131-03	Core	Set functions supported in queries with grouped views	
F131-04	Core	Subqueries with GROUP BY and HAVING clauses and grouped views	
F131-05	Core	Single row SELECT with GROUP BY and HAVING clauses and grouped views	
F171		Multiple schemas per user	
F181	Core	Multiple module support	
F191		Referential delete actions	
F200		TRUNCATE TABLE statement	
F201	Core	CAST function	
F202		TRUNCATE TABLE: identity column restart option	
F221	Core	Explicit defaults	
F222		INSERT statement: DEFAULT VALUES clause	
F231		Privilege tables	
F251		Domain support	
F261	Core	CASE expression	
F261-01	Core	Simple CASE	
F261-02	Core	Searched CASE	
F261-03	Core	NULLIF	
F261-04	Core	COALESCE	

Identifier	Core?	Description	Comment
F262		Extended CASE expression	
F271		Compound character literals	
F281		LIKE enhancements	
F292		UNIQUE null treatment	
F302		INTERSECT table operator	
F303		INTERSECT DISTINCT table operator	
F304		EXCEPT ALL table operator	
F305		INTERSECT ALL table operator	
F311	Core	Schema definition statement	
F311-01	Core	CREATE SCHEMA	
F311-02	Core	CREATE TABLE for persistent base tables	
F311-03	Core	CREATE VIEW	
F311-04	Core	CREATE VIEW: WITH CHECK OPTION	
F311-05	Core	GRANT statement	
F312		MERGE statement	
F313		Enhanced MERGE statement	
F314		MERGE statement with DELETE branch	
F321		User authorization	
F341		Usage tables	
F361		Subprogram support	
F381		Extended schema manipulation	
F382		Alter column data type	
F383		Set column not null clause	
F384		Drop identity property clause	
F385		Drop column generation expression clause	
F386		Set identity column generation clause	
F387		ALTER TABLE statement: ALTER COLUMN clause	
F388		ALTER TABLE statement: ADD/DROP CONSTRAINT clause	
F391		Long identifiers	
F392		Unicode escapes in identifiers	
F393		Unicode escapes in literals	
F394		Optional normal form specification	
F401		Extended joined table	
F402		Named column joins for LOBs, arrays, and multisets	
F404		Range variable for common column names	
F405		NATURAL JOIN	
F406		FULL OUTER JOIN	
F407		CROSS JOIN	
F411		Time zone specification	differences regarding literal interpretation

Identifier	Core?	Description	Comment
F421		National character	
F431		Read-only scrollable cursors	
F432		FETCH with explicit NEXT	
F433		FETCH FIRST	
F434		FETCH LAST	
F435		FETCH PRIOR	
F436		FETCH ABSOLUTE	
F437		FETCH RELATIVE	
F438		Scrollable cursors	
F441		Extended set function support	
F442		Mixed column references in set functions	
F471	Core	Scalar subquery values	
F481	Core	Expanded NULL predicate	
F491		Constraint management	
F501	Core	Features and conformance views	
F501-01	Core	SQL_FEATURES view	
F501-02	Core	SQL_SIZING view	
F502		Enhanced documentation tables	
F531		Temporary tables	
F555		Enhanced seconds precision	
F561		Full value expressions	
F571		Truth value tests	
F591		Derived tables	
F611		Indicator data types	
F641		Row and table constructors	
F651		Catalog name qualifiers	
F661		Simple tables	
F672		Retrospective CHECK constraints	
F690		Collation support	
F692		Extended collation support	
F701		Referential update actions	
F711		ALTER domain	
F731		INSERT column privileges	
F751		View CHECK enhancements	
F761		Session management	
F762		CURRENT_CATALOG	
F763		CURRENT_SCHEMA	
F771		Connection management	
F781		Self-referencing operations	
F791		Insensitive cursors	
F801		Full set function	

Identifier	Core?	Description	Comment
F850		Top-level ORDER BY in query expression	
F851		ORDER BY in subqueries	
F852		Top-level ORDER BY in views	
F855		Nested ORDER BY in query expression	
F856		Nested FETCH FIRST in query expression	
F857		Top-level FETCH FIRST in query expression	
F858		FETCH FIRST in subqueries	
F859		Top-level FETCH FIRST in views	
F860		Dynamic FETCH FIRST row count	
F861		Top-level OFFSET in query expression	
F862		OFFSET in subqueries	
F863		Nested OFFSET in query expression	
F864		Top-level OFFSET in views	
F865		Dynamic offset row count in OFFSET	
F867		FETCH FIRST clause: WITH TIES option	
F868		ORDER BY in grouped table	
F869		SQL implementation info population	
S071		SQL paths in function and type name resolution	
S090		Minimal array support	
S092		Arrays of user-defined types	
S095		Array constructors by query	
S096		Optional array bounds	
S098		ARRAY_AGG	
S099		Array expressions	
S111		ONLY in query expressions	
S201		SQL-invoked routines on arrays	
S203		Array parameters	
S204		Array as result type of functions	
S211		User-defined cast functions	
S301		Enhanced UNNEST	
S404		TRIM_ARRAY	
T031		BOOLEAN data type	
T054		GREATEST and LEAST	different null handling
T055		String padding functions	
T056		Multi-character TRIM functions	
T061		UCS support	
T071		BIGINT data type	
T081		Optional string types maximum length	
T121		WITH (excluding RECURSIVE) in query expression	
T122		WITH (excluding RECURSIVE) in subquery	
T131		Recursive query	

Identifier	Core?	Description	Comment
T132		Recursive query in subquery	
T133		Enhanced cycle mark values	
T141		SIMILAR predicate	
T151		DISTINCT predicate	
T152		DISTINCT predicate with negation	
T171		LIKE clause in table definition	
T172		AS subquery clause in table definition	
T173		Extended LIKE clause in table definition	
T174		Identity columns	
T177		Sequence generator support: simple restart option	
T178		Identity columns: simple restart option	
T191		Referential action RESTRICT	
T201		Comparable data types for referential constraints	
T212		Enhanced trigger capability	
T213		INSTEAD OF triggers	
T214		BEFORE triggers	
T215		AFTER triggers	
T216		Ability to require true search condition before trigger is invoked	
T217		TRIGGER privilege	
T241		START TRANSACTION statement	
T261		Chained transactions	
T271		Savepoints	
T281		SELECT privilege with column granularity	
T285		Enhanced derived column names	
T312		OVERLAY function	
T321-01	Core	User-defined functions with no overloading	
T321-02	Core	User-defined stored procedures with no overloading	
T321-03	Core	Function invocation	
T321-04	Core	CALL statement	
T321-05	Core	RETURN statement	
T321-06	Core	ROUTINES view	
T321-07	Core	PARAMETERS view	
T323		Explicit security for external routines	
T325		Qualified SQL parameter references	
T331		Basic roles	
T332		Extended roles	
T341		Overloading of SQL-invoked functions and SQL-invoked procedures	
T351		Bracketed comments	
T431		Extended grouping capabilities	

Identifier	Core?	Description	Comment
T432		Nested and concatenated GROUPING SETS	
T433		Multi-argument GROUPING function	
T434		GROUP BY DISTINCT	
T441		ABS and MOD functions	
T461		Symmetric BETWEEN predicate	
T491		LATERAL derived table	
T501		Enhanced EXISTS predicate	
T521		Named arguments in CALL statement	
T523		Default values for INOUT parameters of SQL-invoked procedures	
T524		Named arguments in routine invocations other than a CALL statement	
T525		Default values for parameters of SQL-invoked functions	
T551		Optional key words for default syntax	
T581		Regular expression substring function	
T591		UNIQUE constraints of possibly null columns	
T611		Elementary OLAP operations	
T612		Advanced OLAP operations	
T613		Sampling	
T614		NTILE function	
T615		LEAD and LAG functions	
T617		FIRST_VALUE and LAST_VALUE functions	
T620		WINDOW clause: GROUPS option	
T621		Enhanced numeric functions	
T622		Trigonometric functions	
T623		General logarithm functions	
T624		Common logarithm functions	
T626		ANY_VALUE	
T627		Window framed COUNT DISTINCT	
T631	Core	IN predicate with one list element	
T651		SQL-schema statements in SQL routines	
T653		SQL-schema statements in external routines	
T655		Cyclically dependent routines	
T661		Non-decimal integer literals	
T662		Underscores in numeric literals	
T670		Schema and data statement mixing	
T803		String-based JSON	
T811		Basic SQL/JSON constructor functions	
T812		SQL/JSON: JSON_OBJECTAGG	
T813		SQL/JSON: JSON_ARRAYAGG with ORDER BY	
T814		Colon in JSON_OBJECT or JSON_OBJECTAGG	

Identifier	Core?	Description	Comment
T821		Basic SQL/JSON query operators	
T822		SQL/JSON: IS JSON WITH UNIQUE KEYS predicate	
T823		SQL/JSON: PASSING clause	
T824		JSON_TABLE: specific PLAN clause	
T825		SQL/JSON: ON EMPTY and ON ERROR clauses	
T826		General value expression in ON ERROR or ON EMPTY clauses	
T827		JSON_TABLE: sibling NESTED COLUMNS clauses	
T828		JSON_QUERY	
T829		JSON_QUERY: array wrapper options	
T830		Enforcing unique keys in SQL/JSON constructor functions	
T831		SQL/JSON path language: strict mode	
T832		SQL/JSON path language: item method	
T833		SQL/JSON path language: multiple subscripts	
T834		SQL/JSON path language: wildcard member accessor	
T835		SQL/JSON path language: filter expressions	
T836		SQL/JSON path language: starts with predicate	
T837		SQL/JSON path language: regex_like predicate	
T838		JSON_TABLE: PLAN DEFAULT clause	
T840		Hex integer literals in SQL/JSON path language	
T851		SQL/JSON: optional keywords for default syntax	
T865		SQL/JSON item method: bigint()	
T866		SQL/JSON item method: boolean()	
T867		SQL/JSON item method: date()	
T868		SQL/JSON item method: decimal()	
T869		SQL/JSON item method: decimal() with precision and scale	
T870		SQL/JSON item method: integer()	
T871		SQL/JSON item method: number()	
T872		SQL/JSON item method: string()	
T873		SQL/JSON item method: time()	
T874		SQL/JSON item method: time_tz()	
T875		SQL/JSON item method: time precision	
T876		SQL/JSON item method: timestamp()	
T877		SQL/JSON item method: timestamp_tz()	
T878		SQL/JSON item method: timestamp precision	
T879		JSON in equality operations	with jsonb
T880		JSON in grouping operations	with jsonb
X010		XML type	
X011		Arrays of XML type	
X014		Attributes of XML type	

Identifier	Core?	Description	Comment
X016		Persistent XML values	
X020		XMLConcat	
X031		XMLElement	
X032		XMLForest	
X034		XMLAgg	
X035		XMLAgg: ORDER BY option	
X036		XMLComment	
X037		XMLPI	
X038		XMLText	supported except for RETURNING
X040		Basic table mapping	
X041		Basic table mapping: null absent	
X042		Basic table mapping: null as nil	
X043		Basic table mapping: table as forest	
X044		Basic table mapping: table as element	
X045		Basic table mapping: with target namespace	
X046		Basic table mapping: data mapping	
X047		Basic table mapping: metadata mapping	
X048		Basic table mapping: base64 encoding of binary strings	
X049		Basic table mapping: hex encoding of binary strings	
X050		Advanced table mapping	
X051		Advanced table mapping: null absent	
X052		Advanced table mapping: null as nil	
X053		Advanced table mapping: table as forest	
X054		Advanced table mapping: table as element	
X055		Advanced table mapping: with target namespace	
X056		Advanced table mapping: data mapping	
X057		Advanced table mapping: metadata mapping	
X058		Advanced table mapping: base64 encoding of binary strings	
X059		Advanced table mapping: hex encoding of binary strings	
X060		XMLParse: character string input and CONTENT option	
X061		XMLParse: character string input and DOCUMENT option	
X069		XMLSerialize: INDENT	
X070		XMLSerialize: character string serialization and CONTENT option	
X071		XMLSerialize: character string serialization and DOCUMENT option	
X072		XMLSerialize: character string serialization	
X090		XML document predicate	
X120		XML parameters in SQL routines	

Identifier	Core?	Description	Comment
X121		XML parameters in external routines	
X221		XML passing mechanism BY VALUE	
X301		XMLTable: derived column list option	
X302		XMLTable: ordinality column option	
X303		XMLTable: column default option	
X304		XMLTable: passing a context item	must be XML DOCUMENT
X400		Name and identifier mapping	
X410		Alter column data type: XML type	

D.2. Unsupported Features

The following features defined in SQL:2023 are not implemented in this release of Postgres Pro. In a few cases, equivalent functionality is available.

Identifier	Core?	Description	Comment
B011		Embedded Ada	
B013		Embedded COBOL	
B014		Embedded Fortran	
B015		Embedded MUMPS	
B016		Embedded Pascal	
B017		Embedded PL/I	
B030		Enhanced dynamic SQL	
B031		Basic dynamic SQL	
B032		Extended dynamic SQL	
B033		Untyped SQL-invoked function arguments	
B034		Dynamic specification of cursor attributes	
B035		Non-extended descriptor names	
B036		Describe input statement	
B041		Extensions to embedded SQL exception declarations	
B051		Enhanced execution rights	
B111		Module language Ada	
B112		Module language C	
B113		Module language COBOL	
B114		Module language Fortran	
B115		Module language MUMPS	
B116		Module language Pascal	
B117		Module language PL/I	
B121		Routine language Ada	
B122		Routine language C	
B123		Routine language COBOL	
B124		Routine language Fortran	
B125		Routine language MUMPS	

Identifier	Core?	Description	Comment
B126		Routine language Pascal	
B127		Routine language PL/I	
B200		Polymorphic table functions	
B201		More than one PTF generic table parameter	
B202		PTF copartitioning	
B203		More than one copartition specification	
B204		PRUNE WHEN EMPTY	
B205		Pass-through columns	
B206		PTF descriptor parameters	
B207		Cross products of partitionings	
B208		PTF component procedure interface	
B209		PTF extended names	
B211		Module language Ada: VARCHAR and NUMERIC support	
B221		Routine language Ada: VARCHAR and NUMERIC support	
F054		TIMESTAMP in DATE type precedence list	
F120		Get diagnostics statement	
F121		Basic diagnostics management	
F122		Enhanced diagnostics management	
F123		All diagnostics	
F124		SET TRANSACTION statement: DIAGNOSTICS SIZE clause	
F263		Comma-separated predicates in simple CASE expression	
F291		UNIQUE predicate	
F301		CORRESPONDING in query expressions	
F403		Partitioned join tables	
F451		Character set definition	
F461		Named character sets	
F492		Optional table constraint enforcement	
F521		Assertions	
F671		Subqueries in CHECK constraints	intentionally omitted
F673		Reads SQL-data routine invocations in CHECK constraints	
F693		SQL-session and client module collations	
F695		Translation support	
F696		Additional translation documentation	
F721		Deferrable constraints	foreign and unique keys only
F741		Referential MATCH types	no partial match yet
F812		Basic flagging	
F813		Extended flagging	

Identifier	Core?	Description	Comment
F821		Local table references	
F831		Full cursor update	
F832		Updatable scrollable cursors	
F833		Updatable ordered cursors	
F841		LIKE_REGEX predicate	consider regexp_like()
F842		OCCURRENCES_REGEX function	consider regexp_matches()
F843		POSITION_REGEX function	consider regexp_instr()
F844		SUBSTRING_REGEX function	consider regexp_substr()
F845		TRANSLATE_REGEX function	consider regexp_replace()
F846		Octet support in regular expression operators	
F847		Non-constant regular expressions	
F866		FETCH FIRST clause: PERCENT option	
R010		Row pattern recognition: FROM clause	
R020		Row pattern recognition: WINDOW clause	
R030		Row pattern recognition: full aggregate support	
S011	Core	Distinct data types	
S011-01	Core	USER_DEFINED_TYPES view	
S023		Basic structured types	
S024		Enhanced structured types	
S025		Final structured types	
S026		Self-referencing structured types	
S027		Create method by specific method name	
S028		Permutable UDT options list	
S041		Basic reference types	
S043		Enhanced reference types	
S051		Create table of type	partially supported
S081		Subtables	
S091		Basic array support	partially supported
S093		Arrays of distinct types	
S094		Arrays of reference types	
S097		Array element assignment	
S151		Type predicate	see pg_typeof()
S161		Subtype treatment	
S162		Subtype treatment for references	
S202		SQL-invoked routines on multisets	
S231		Structured type locators	
S232		Array locators	
S233		Multiset locators	
S241		Transform functions	

Identifier	Core?	Description	Comment
S242		Alter transform statement	
S251		User-defined orderings	
S261		Specific type method	
S271		Basic multiset support	
S272		Multisets of user-defined types	
S274		Multisets of reference types	
S275		Advanced multiset support	
S281		Nested collection types	
S291		Unique constraint on entire row	
S401		Distinct types based on array types	
S402		Distinct types based on multiset types	
S403		ARRAY_MAX_CARDINALITY	
T011		Timestamp in Information Schema	
T021		BINARY and VARBINARY data types	
T022		Advanced support for BINARY and VARBINARY data types	
T023		Compound binary literals	
T024		Spaces in binary literals	
T039		CLOB locator: non-holdable	
T040		Concatenation of CLOBs	
T041		Basic LOB data type support	
T042		Extended LOB data type support	
T043		Multiplier T	
T044		Multiplier P	
T045		BLOB data type	
T046		CLOB data type	
T047		POSITION, OCTET_LENGTH, TRIM, and SUBSTRING for BLOBs	
T048		Concatenation of BLOBs	
T049		BLOB locator: non-holdable	
T050		POSITION, CHAR_LENGTH, OCTET_LENGTH, LOWER, TRIM, UPPER, and SUBSTRING for CLOBs	
T051		Row types	
T053		Explicit aliases for all-fields reference	
T062		Character length units	
T076		DECFLOAT data type	
T101		Enhanced nullability determination	
T111		Updatable joins, unions, and columns	
T175		Generated columns	mostly supported
T176		Sequence generator support	supported except for NEXT VALUE FOR
T180		System-versioned tables	

Identifier	Core?	Description	Comment
T181		Application-time period tables	
T200		Trigger DDL	similar but not fully compatible
T211		Basic trigger capability	
T218		Multiple triggers for the same event executed in the order created	intentionally omitted
T231		Sensitive cursors	
T251		SET TRANSACTION statement: LOCAL option	
T262		Multiple server transactions	
T272		Enhanced savepoint management	
T301		Functional dependencies	partially supported
T321	Core	Basic SQL-invoked routines	partially supported
T322		Declared data type attributes	
T324		Explicit security for SQL routines	
T326		Table functions	
T471		Result sets return value	
T472		DESCRIBE CURSOR	
T495		Combined data change and retrieval	different syntax
T502		Period predicates	
T511		Transaction counts	
T522		Default values for IN parameters of SQL-invoked procedures	supported except DEFAULT key word in invocation
T561		Holdable locators	
T571		Array-returning external SQL-invoked functions	
T572		Multiset-returning external SQL-invoked functions	
T601		Local cursor references	
T616		Null treatment option for LEAD and LAG functions	
T618		NTH_VALUE function	function exists, but some options missing
T619		Nested window functions	
T625		LISTAGG	
T641		Multiple column assignment	only some syntax variants supported
T652		SQL-dynamic statements in SQL routines	
T654		SQL-dynamic statements in external routines	
T801		JSON data type	
T802		Enhanced JSON data type	
T839		Formatted cast of datetimes to/from character strings	
T860		SQL/JSON simplified accessor: column reference only	
T861		SQL/JSON simplified accessor: case-sensitive JSON member accessor	

Identifier	Core?	Description	Comment
T862		SQL/JSON simplified accessor: wildcard member accessor	
T863		SQL/JSON simplified accessor: single-quoted string literal as member accessor	
T864		SQL/JSON simplified accessor	
T881		JSON in ordering operations	with jsonb, partially supported
T882		JSON in multiset element grouping operations	
M001		Datalinks	
M002		Datalinks via SQL/CLI	
M003		Datalinks via Embedded SQL	
M004		Foreign data support	partially supported
M005		Foreign schema support	
M006		GetSQLString routine	
M007		TransmitRequest	
M009		GetOpts and GetStatistics routines	
M010		Foreign-data wrapper support	different API
M011		Datalinks via Ada	
M012		Datalinks via C	
M013		Datalinks via COBOL	
M014		Datalinks via Fortran	
M015		Datalinks via M	
M016		Datalinks via Pascal	
M017		Datalinks via PL/I	
M018		Foreign-data wrapper interface routines in Ada	
M019		Foreign-data wrapper interface routines in C	different API
M020		Foreign-data wrapper interface routines in COBOL	
M021		Foreign-data wrapper interface routines in Fortran	
M022		Foreign-data wrapper interface routines in MUMPS	
M023		Foreign-data wrapper interface routines in Pascal	
M024		Foreign-data wrapper interface routines in PL/I	
M030		SQL-server foreign data support	
M031		Foreign-data wrapper general routines	
X012		Multisets of XML type	
X013		Distinct types of XML type	
X015		Fields of XML type	
X025		XMLCast	
X030		XMLDocument	
X065		XMLParse: binary string input and CONTENT option	
X066		XMLParse: binary string input and DOCUMENT option	
X068		XMLSerialize: BOM	

Identifier	Core?	Description	Comment
X073		XMLSerialize: binary string serialization and CONTENT option	
X074		XMLSerialize: binary string serialization and DOCUMENT option	
X075		XMLSerialize: binary string serialization	
X076		XMLSerialize: VERSION	
X077		XMLSerialize: explicit ENCODING option	
X078		XMLSerialize: explicit XML declaration	
X080		Namespaces in XML publishing	
X081		Query-level XML namespace declarations	
X082		XML namespace declarations in DML	
X083		XML namespace declarations in DDL	
X084		XML namespace declarations in compound statements	
X085		Predefined namespace prefixes	
X086		XML namespace declarations in XMLTable	
X091		XML content predicate	
X096		XMLExists	XPath 1.0 only
X100		Host language support for XML: CONTENT option	
X101		Host language support for XML: DOCUMENT option	
X110		Host language support for XML: VARCHAR mapping	
X111		Host language support for XML: CLOB mapping	
X112		Host language support for XML: BLOB mapping	
X113		Host language support for XML: STRIP WHITESPACE option	
X114		Host language support for XML: PRESERVE WHITESPACE option	
X131		Query-level XMLBINARY clause	
X132		XMLBINARY clause in DML	
X133		XMLBINARY clause in DDL	
X134		XMLBINARY clause in compound statements	
X135		XMLBINARY clause in subqueries	
X141		IS VALID predicate: data-driven case	
X142		IS VALID predicate: ACCORDING TO clause	
X143		IS VALID predicate: ELEMENT clause	
X144		IS VALID predicate: schema location	
X145		IS VALID predicate outside check constraints	
X151		IS VALID predicate: with DOCUMENT option	
X152		IS VALID predicate: with CONTENT option	
X153		IS VALID predicate: with SEQUENCE option	
X155		IS VALID predicate: NAMESPACE without ELEMENT clause	

Identifier	Core?	Description	Comment
X157		IS VALID predicate: NO NAMESPACE with ELEMENT clause	
X160		Basic Information Schema for registered XML schemas	
X161		Advanced Information Schema for registered XML schemas	
X170		XML null handling options	
X171		NIL ON NO CONTENT option	
X181		XML(DOCUMENT(UNTYPED)) type	
X182		XML(DOCUMENT(ANY)) type	
X190		XML(SEQUENCE) type	
X191		XML(DOCUMENT(XMLSCHEMA)) type	
X192		XML(CONTENT(XMLSCHEMA)) type	
X200		XMLQuery	
X201		XMLQuery: RETURNING CONTENT	
X202		XMLQuery: RETURNING SEQUENCE	
X203		XMLQuery: passing a context item	
X204		XMLQuery: initializing an XQuery variable	
X205		XMLQuery: EMPTY ON EMPTY option	
X206		XMLQuery: NULL ON EMPTY option	
X211		XML 1.1 support	
X222		XML passing mechanism BY REF	parser accepts BY REF but ignores it; passing is always BY VALUE
X231		XML(CONTENT(UNTYPED)) type	
X232		XML(CONTENT(ANY)) type	
X241		RETURNING CONTENT in XML publishing	
X242		RETURNING SEQUENCE in XML publishing	
X251		Persistent XML values of XML(DOCUMENT(UNTYPED)) type	
X252		Persistent XML values of XML(DOCUMENT(ANY)) type	
X253		Persistent XML values of XML(CONTENT(UNTYPED)) type	
X254		Persistent XML values of XML(CONTENT(ANY)) type	
X255		Persistent XML values of XML(SEQUENCE) type	
X256		Persistent XML values of XML(DOCUMENT(XMLSCHEMA)) type	
X257		Persistent XML values of XML(CONTENT(XMLSCHEMA)) type	
X260		XML type: ELEMENT clause	
X261		XML type: NAMESPACE without ELEMENT clause	
X263		XML type: NO NAMESPACE with ELEMENT clause	
X264		XML type: schema location	
X271		XMLValidate: data-driven case	

Identifier	Core?	Description	Comment
X272		XMLValidate: ACCORDING TO clause	
X273		XMLValidate: ELEMENT clause	
X274		XMLValidate: schema location	
X281		XMLValidate with DOCUMENT option	
X282		XMLValidate with CONTENT option	
X283		XMLValidate with SEQUENCE option	
X284		XMLValidate: NAMESPACE without ELEMENT clause	
X286		XMLValidate: NO NAMESPACE with ELEMENT clause	
X300		XMLTable	XPath 1.0 only
X305		XMLTable: initializing an XQuery variable	

D.3. XML Limits and Conformance to SQL/XML

Significant revisions to the XML-related specifications in ISO/IEC 9075-14 (SQL/XML) were introduced with SQL:2006. Postgres Pro's implementation of the XML data type and related functions largely follows the earlier 2003 edition, with some borrowing from later editions. In particular:

- Where the current standard provides a family of XML data types to hold “document” or “content” in untyped or XML Schema-typed variants, and a type `XML(SEQUENCE)` to hold arbitrary pieces of XML content, Postgres Pro provides the single `xml` type, which can hold “document” or “content”. There is no equivalent of the standard's “sequence” type.
- Postgres Pro provides two functions introduced in SQL:2006, but in variants that use the XPath 1.0 language, rather than XML Query as specified for them in the standard.
- Postgres Pro does not support the `RETURNING CONTENT` or `RETURNING SEQUENCE` clauses, functions which are defined to have these in the specification are implicitly returning content.

This section presents some of the resulting differences you may encounter.

D.3.1. Queries Are Restricted to XPath 1.0

The Postgres Pro-specific functions `xpath()` and `xpath_exists()` query XML documents using the XPath language. Postgres Pro also provides XPath-only variants of the standard functions `XMLEXISTS` and `XMLTABLE`, which officially use the XQuery language. For all of these functions, Postgres Pro relies on the libxml2 library, which provides only XPath 1.0.

There is a strong connection between the XQuery language and XPath versions 2.0 and later: any expression that is syntactically valid and executes successfully in both produces the same result (with a minor exception for expressions containing numeric character references or predefined entity references, which XQuery replaces with the corresponding character while XPath leaves them alone). But there is no such connection between these languages and XPath 1.0; it was an earlier language and differs in many respects.

There are two categories of limitation to keep in mind: the restriction from XQuery to XPath for the functions specified in the SQL standard, and the restriction of XPath to version 1.0 for both the standard and the Postgres Pro-specific functions.

D.3.1.1. Restriction of XQuery to XPath

Features of XQuery beyond those of XPath include:

- XQuery expressions can construct and return new XML nodes, in addition to all possible XPath values. XPath can create and return values of the atomic types (numbers, strings, and so on) but can only return XML nodes that were already present in documents supplied as input to the expression.
- XQuery has control constructs for iteration, sorting, and grouping.

- XQuery allows declaration and use of local functions.

Recent XPath versions begin to offer capabilities overlapping with these (such as functional-style `for-each` and `sort`, anonymous functions, and `parse-xml` to create a node from a string), but such features were not available before XPath 3.0.

D.3.1.2. Restriction of XPath to 1.0

For developers familiar with XQuery and XPath 2.0 or later, XPath 1.0 presents a number of differences to contend with:

- The fundamental type of an XQuery/XPath expression, the *sequence*, which can contain XML nodes, atomic values, or both, does not exist in XPath 1.0. A 1.0 expression can only produce a node-set (containing zero or more XML nodes), or a single atomic value.
- Unlike an XQuery/XPath sequence, which can contain any desired items in any desired order, an XPath 1.0 node-set has no guaranteed order and, like any set, does not allow multiple appearances of the same item.

Note

The libxml2 library does seem to always return node-sets to Postgres Pro with their members in the same relative order they had in the input document. Its documentation does not commit to this behavior, and an XPath 1.0 expression cannot control it.

- While XQuery/XPath provides all of the types defined in XML Schema and many operators and functions over those types, XPath 1.0 has only node-sets and the three atomic types `boolean`, `double`, and `string`.
- XPath 1.0 has no conditional operator. An XQuery/XPath expression such as `if (hat) then hat/@size else "no hat"` has no XPath 1.0 equivalent.
- XPath 1.0 has no ordering comparison operator for strings. Both `"cat" < "dog"` and `"cat" > "dog"` are false, because each is a numeric comparison of two NaNs. In contrast, `=` and `!=` do compare the strings as strings.
- XPath 1.0 blurs the distinction between *value comparisons* and *general comparisons* as XQuery/XPath define them. Both `sale/@hatsize = 7` and `sale/@customer = "alice"` are existentially quantified comparisons, true if there is any `sale` with the given value for the attribute, but `sale/@taxable = false()` is a value comparison to the *effective boolean value* of a whole node-set. It is true only if no `sale` has a `taxable` attribute at all.
- In the XQuery/XPath data model, a *document node* can have either document form (i.e., exactly one top-level element, with only comments and processing instructions outside of it) or content form (with those constraints relaxed). Its equivalent in XPath 1.0, the *root node*, can only be in document form. This is part of the reason an `xml` value passed as the context item to any Postgres Pro XPath-based function must be in document form.

The differences highlighted here are not all of them. In XQuery and the 2.0 and later versions of XPath, there is an XPath 1.0 compatibility mode, and the W3C lists of [function library changes](#) and [language changes](#) applied in that mode offer a more complete (but still not exhaustive) account of the differences. The compatibility mode cannot make the later languages exactly equivalent to XPath 1.0.

D.3.1.3. Mappings between SQL and XML Data Types and Values

In SQL:2006 and later, both directions of conversion between standard SQL data types and the XML Schema types are specified precisely. However, the rules are expressed using the types and semantics of XQuery/XPath, and have no direct application to the different data model of XPath 1.0.

When Postgres Pro maps SQL data values to XML (as in `xmlelement`), or XML to SQL (as in the output columns of `xmltable`), except for a few cases treated specially, Postgres Pro simply assumes that the

XML data type's XPath 1.0 string form will be valid as the text-input form of the SQL datatype, and conversely. This rule has the virtue of simplicity while producing, for many data types, results similar to the mappings specified in the standard.

Where interoperability with other systems is a concern, for some data types, it may be necessary to use data type formatting functions (such as those in [Section 9.8](#)) explicitly to produce the standard mappings.

D.3.2. Incidental Limits of the Implementation

This section concerns limits that are not inherent in the libxml2 library, but apply to the current implementation in Postgres Pro.

D.3.2.1. Only `BY VALUE` Passing Mechanism Is Supported

The SQL standard defines two *passing mechanisms* that apply when passing an XML argument from SQL to an XML function or receiving a result: `BY REF`, in which a particular XML value retains its node identity, and `BY VALUE`, in which the content of the XML is passed but node identity is not preserved. A mechanism can be specified before a list of parameters, as the default mechanism for all of them, or after any parameter, to override the default.

To illustrate the difference, if `x` is an XML value, these two queries in an SQL:2006 environment would produce true and false, respectively:

```
SELECT XMLQUERY('$a is $b' PASSING BY REF x AS a, x AS b NULL ON EMPTY);
SELECT XMLQUERY('$a is $b' PASSING BY VALUE x AS a, x AS b NULL ON EMPTY);
```

Postgres Pro will accept `BY VALUE` or `BY REF` in an `XML EXISTS` or `XMLTABLE` construct, but it ignores them. The `xml` data type holds a character-string serialized representation, so there is no node identity to preserve, and passing is always effectively `BY VALUE`.

D.3.2.2. Cannot Pass Named Parameters to Queries

The XPath-based functions support passing one parameter to serve as the XPath expression's context item, but do not support passing additional values to be available to the expression as named parameters.

D.3.2.3. No `XML (SEQUENCE)` Type

The Postgres Pro `xml` data type can only hold a value in `DOCUMENT` or `CONTENT` form. An XQuery/XPath expression context item must be a single XML node or atomic value, but XPath 1.0 further restricts it to be only an XML node, and has no node type allowing `CONTENT`. The upshot is that a well-formed `DOCUMENT` is the only form of XML value that Postgres Pro can supply as an XPath context item.

Appendix E. Release Notes

The release notes contain the significant changes in each Postgres Pro Shardman release, with major features and migration issues listed at the top. The release notes do not contain changes that affect only a few users or changes that are internal and therefore not user-visible. For example, the optimizer is improved in almost every release, but the improvements are usually observed by users as simply faster queries.

E.1. Postgres Pro Shardman 17.5.1

Release date: 2025-07-02

This is the first Postgres Pro Shardman release, it is based on Postgres Pro Enterprise 17.5.1 and Shardman 14.18.2 and provides new features and optimizations. Major changes are as follows:

E.1.1. Core and Extensions

- Implemented the `LOCK TABLE` support for global and sharded tables.
- Renamed the `is_reply` field of the `shardman.silk_routing` view to `is_outgoing`. The field value is `true` if the route was produced by the outgoing network connection, `false` for incoming network connections.
- Enhanced Silk performance for higher `shardman.silk_unassigned_job_queue_size` values.
- Updated the default values for logging configuration parameters for all new Postgres Pro Shardman installations:
 - `logging_collector`: on
 - `log_directory`: `{{dataDir}}`
 - `log_filename`: `postgres-{{cluster}}-{{shard}}-{{host}}-{{keeperID}}.log`
- Deleted the `csn_max_commit_shift` and `csn_max_snapshot_shift` configuration parameters. Added the `csn_max_shift` and `csn_max_shift_error` configuration parameters to work with CSN snapshots for the distributed queries and imported snapshots.
- Added the `background_freezer` feature that vacuums modified and non-locked data pages of permanent tables stored in RAM. Background freezer is managed using new specific [configuration parameters](#).
- Changed lock logic to reduce CPU consumption significantly in case of startup recovery on physical replicas used as a source for logical replication with many downstreams. Walsenders no longer slow down WAL application as a more efficient synchronization primitive has replaced the previously used spinlock.
- Added the new `planner_upper_limit_estimation` configuration parameter that enables the query planner to overestimate the expected number of rows in statements that contain a comparison to an unknown constant.
- Added a possibility for the autovacuum daemon to process indexes of a table in a parallel mode. The `parallel_autovacuum_workers` storage parameter of a table controls whether it should be processed in parallel and defines the number of additional parallel autovacuum workers that can be launched for table processing. The number of autovacuum workers per cluster that can be used for parallel table processing is limited by the `autovacuum_max_workers` configuration parameter. Parallel autovacuum is currently in an experimental phase and is intended for testing purposes only. Do not deploy it in production environments.
- Added the `enable_any_to_lateral_transformation` configuration parameter that allows enabling or disabling transformation of `ANY` subqueries into `LATERAL` joins.
- Implemented performance optimization when working with table metadata, which allows obtaining information about attributes using the system cache instead of direct reading from the system catalog.

- Introduced the following changes to the implementation of `crash_info`:
 - Added handling of `SIGILL` signals to `crash_info` processing.
 - Added details, like process start time and query text at planning, to `crash_info` output files.
- Restricted superuser actions with temporary relations of other sessions. Superusers can now only use the `DROP TABLE` command with such relations. If the `skip_temp_rel_lock` configuration parameter is set to `on`, even dropping such relations is not allowed.
- Upgraded `pg_proaudit` to provide the following enhancements:
 - Added new object types: `CATALOG RELATION` and `CATALOG FUNCTION`.
 - Added new event fields: `UUID`, `XID`, and `VXID`. Now it is possible to identify an event by its `UUID` and transaction IDs (if applicable).
 - The behavior of the `pg_proaudit.log_catalog_access` configuration parameter was corrected to reflect the new logic of logging for system catalog objects.
 - Changed the logic of handling disconnection events. Now these events are associated with the corresponding authentication events, so the disconnection events will be logged even in the case when the rule is removed after authentication, but before disconnect.
- Upgraded `pgpro_pwr` to version 4.9, which mainly provides optimizations. Notable changes are as follows:
 - Added support for `pgpro_stats 1.9-sdm4`.
 - Added a possibility to define the mode of collecting relation sizes globally through the `pgpro_pwr.relsizes_collect_mode` extension parameter or for a server through the `set_server_size_sampling` function parameter.
 - Enabled fine-tuning the server statistics collection by calling the `set_server_setting` function. It allows you to define which statistics should be collected.
 - Added a preview of storage parameters for tables and indexes in the “Schema object statistics” report section.
- Upgraded `pgpro_stats` to version 1.9-sdm4, which provides the following improvements:
 - Enhanced session tracing to provide more information. Specifically, the new `time_info` filter attribute controls inclusion of additional information in the session-tracing output, and the `pgpro_stats.trace_query_text_size` configuration parameter can limit the size of the query in the session-tracing output.
 - Aligned the names of the `explain_*` filter attributes of the session tracer with the names of session-tracing configuration parameters.
 - Changed the format of the statistics dump file and the corresponding save/load routines.
 - Implemented turning off the session tracer functionality when no session-tracing filters are specified.
 - Prohibited inclusion of both `pgpro_stats` and `pg_stat_statements` in the list of `shared_preload_libraries`. If both are included, the database server will not start.
- Deprecated the experimental `vops` extension.

E.1.2. Management Utilities

- Changed the default values for the `-t` (`--transactions`) and `-T` (`--time`) flags of the `shardmanctl bench run` command to 10 and 0, respectively.
- Optimized promotion of a standby server to primary.
- Added a limitation that `pg_hba.conf` cannot be modified if `certVersion` is set to `true`.

- Added the `--reinit-before-restore` option that allows reinitializing the cluster with the same topology before restore.
- Optimized cluster polling timeouts.
- Added the `{rgid}` support for the `--tablespace-mapping` of the `shardmanctl probackup restore` command.
- Added the feature to validate a backup on a compatible cluster.
- Optimized the search for the primary node to avoid cluster hanging.
- Added `boot_uuid` to `shardmand`.
- Added the feature to modify the `max_wal_senders` parameter by updating the `additionalWalSenders` configuration parameter with the `shardmanctl config update` command.
- Added the `shardmanctl --log-format` command to specify the log output format.
- Updated the `shardmanctl cluster stop` and `shardmanctl restart` commands with the `--message` option to show the reason why a cluster was stopped or restarted.
- Added the `CertVersion` field to the `sdmspec.json` configuration file.
- Updated the `syncPointMonitorEnabled` configuration parameter to save the syncpoint history if set to `true`, with size of this history defined by `syncPointMonitorLimit`.
- Updated the `keeper` error messages that are now cleaned after a successful operation that modifies the configuration.
- Added the `shardmanctl schema verify` command that verifies consistency of the information about partitions, extensions, tables, roles, schemas, and sequences.
- Added an option to use `probackup` as backup method for BiHA.
- Added a feature that cluster software does not fail even if writing to `pg_hba.conf` is forbidden.

E.2. Release 17.5

Release date: 2025-05-08

This release contains a variety of fixes from 17.4. For information about new features in major release 17, see [Section E.7](#).

E.2.1. Migration to Version 17.5

A dump/restore is not required for those running 17.X.

However, if you have any self-referential foreign key constraints on partitioned tables, it may be necessary to recreate those constraints to ensure that they are being enforced correctly. See the second changelog entry below.

Also, if you have any BRIN bloom indexes, it may be advisable to reindex them after updating. See the third changelog entry below.

Also, if you are upgrading from a version earlier than 17.1, see [Section E.6](#).

E.2.2. Changes

- Avoid one-byte buffer overread when examining invalidly-encoded strings that are claimed to be in GB18030 encoding (Noah Misch, Andres Freund)

While unlikely, a SIGSEGV crash could occur if an incomplete multibyte character appeared at the end of memory. This was possible both in the server and in libpq-using applications. (CVE-2025-4207)

- Handle self-referential foreign keys on partitioned tables correctly (Álvaro Herrera)

Creating or attaching partitions failed to make the required catalog entries for a foreign-key constraint, if the table referenced by the constraint was the same partitioned table. This resulted in failure to enforce the constraint fully.

To fix this, you should drop and recreate any self-referential foreign keys on partitioned tables, if partitions have been created or attached since the constraint was created. Bear in mind that violating rows might already be present, in which case recreating the constraint will fail, and you'll need to fix up those rows before trying again.

- Avoid data loss when merging compressed BRIN summaries in `brin_bloom_union()` (Tomas Vondra)

The code failed to account for decompression results not being identical to the input objects, which would result in failure to add some of the data to the merged summary, leading to missed rows in index searches.

This mistake was present back to v14 where BRIN bloom indexes were introduced, but this code path was only rarely reached then. It's substantially more likely to be hit in v17 because parallel index builds now use the code.

- Fix unexpected “attribute has wrong type” errors in `UPDATE`, `DELETE`, and `MERGE` queries that use whole-row table references to views or functions in `FROM` (Tom Lane)
- Fix `MERGE` into a partitioned table with `DO NOTHING` actions (Tender Wang)

Some cases failed with “unknown action in `MERGE WHEN` clause” errors.

- Prevent failure in `INSERT` commands when the table has a `GENERATED` column of a domain data type and the domain's constraints disallow null values (Jian He)

Constraint failure was reported even if the generation expression produced a perfectly okay result.

- Correctly process references to outer CTE names that appear within a `WITH` clause attached to an `INSERT/UPDATE/DELETE/MERGE` command that's inside `WITH` (Tom Lane)

The parser failed to detect disallowed recursion cases, nor did it account for such references when sorting CTEs into a usable order.

- Fix misprocessing of casts within the keys of JSON constructor expressions (Amit Langote)
- Don't try to parallelize `array_agg()` when the argument is of an anonymous record type (Richard Guo, Tom Lane)

The protocol for communicating with parallel workers doesn't support identifying the concrete record type that a worker is returning.

- Fix `ARRAY(subquery)` and `ARRAY[expression, ...]` constructs to produce sane results when the input is of type `int2vector` or `oidvector` (Tom Lane)

This patch restores the behavior that existed before PostgreSQL 9.5: the result is of type `int2vector[]` or `oidvector[]`.

- Fix possible erroneous reports of invalid affixes while parsing Ispell dictionaries (Jacob Brazeal)
- Fix `ALTER TABLE ADD COLUMN` to correctly handle the case of a domain type that has a default (Jian He, Tom Lane, Tender Wang)

If a domain type has a default, adding a column of that type (without any explicit `DEFAULT` clause) failed to install the domain's default value in existing rows, instead leaving the new column null.

- Repair misbehavior when there are duplicate column names in a foreign key constraint's `ON DELETE SET DEFAULT` or `SET NULL` action (Tom Lane)

- Improve the error message for disallowed attempts to alter the properties of a foreign key constraint (Álvaro Herrera)
- Avoid error when resetting the `relhassubclass` flag of a temporary table that's marked `ON COMMIT DELETE ROWS` (Noah Misch)
- Add missing deparsing of the `INDENT` option of `XMLSERIALIZE()` (Jim Jones)

Previously, views or rules using `XMLSERIALIZE(... INDENT)` were dumped without the `INDENT` clause, causing incorrect results after restore.

- Avoid premature evaluation of the arguments of an aggregate function that has both `FILTER` and `ORDER BY` (or `DISTINCT`) options (David Rowley)

If there is `ORDER BY` or `DISTINCT`, we consider pre-sorting the aggregate input values rather than doing the sort within the Agg plan node. But this is problematic if the aggregate inputs include expressions that could fail (for example, a division where some of the input divisors could be zero) and there is a `FILTER` clause that's meant to prevent such failures. Pre-sorting would push the expression evaluations to before the `FILTER` test, allowing the failures to happen anyway. Avoid this by not pre-sorting if there's a `FILTER` and the input expressions are anything more complex than a simple Var or Const.

- Fix erroneous deductions from column `NOT NULL` constraints in the presence of outer joins (Richard Guo)

In some cases the planner would discard an `IS NOT NULL` query condition, even though the condition applies after an outer join and thus is not redundant.

- Avoid incorrect optimizations based on `IS [NOT] NULL` tests that are applied to composite values (Bruce Momjian)
- Fix planner's failure to identify more than one hashable `ScalarArrayOpExpr` subexpression within a top-level expression (David Geier)

This resulted in unnecessarily-inefficient execution of any additional subexpressions that could have been processed with a hash table (that is, `IN`, `NOT IN`, or `= ANY` clauses with all-constant right-hand sides).

- Fix incorrect table size estimate with low fill factor (Tomas Vondra)

When the planner estimates the number of rows in a never-yet-analyzed table, it uses the table's `fillfactor` setting in the estimation, but it neglected to clamp the result to at least one row per page. A low `fillfactor` could thus result in an unreasonably small estimate.

- Disable “skip fetch” optimization in bitmap heap scan (Matthias van de Meent)

It turns out that this optimization can result in returning dead tuples when a concurrent vacuum marks a page all-visible.

- Fix performance issues in GIN index search startup when there are many search keys (Tom Lane, Vinod Sridharan)

An indexable clause with many keys (for example, `jsonbcol ?| array[...]` with tens of thousands of array elements) took $O(N^2)$ time to start up, and was uncancelable for that interval too.

- Detect missing support procedures in a BRIN index operator class, and report an error instead of crashing (Álvaro Herrera)
- Respond to interrupts (such as query cancel) while waiting for asynchronous subplans of an Append plan node (Heikki Linnakangas)

Previously, nothing would happen until one of the subplans becomes ready.

- Report the I/O statistics of active WAL senders more frequently (Bertrand Drouvot)

Previously, the `pg_stat_io` view failed to accumulate I/O performed by a WAL sender until that process exited. Now such I/O will be reported after at most one second's delay.

- Fix race condition in handling of `synchronous_standby_names` immediately after startup (Melnikov Maksim, Michael Paquier)

For a short period after system startup, backends might fail to wait for synchronous commit even though `synchronous_standby_names` is enabled.

- Cope with possible intra-query changes of `io_combine_limit` (Thomas Munro)
- Avoid infinite loop if `scram_iterations` is set to `INT_MAX` (Kevin K Bijou)
- Avoid possible crashes due to double transformation of `json_array()`'s subquery (Tom Lane)
- Fix `pg_strtof()` to not crash with null `endptr` (Alexander Lakhin, Tom Lane)
- Fix crash after out-of-memory in certain GUC assignments (Daniel Gustafsson)
- Avoid crash when a Snowball stemmer encounters an out-of-memory condition (Maksim Korotkov)
- Fix over-enthusiastic freeing of `SpecialJoinInfo` structs during planning (Richard Guo)

This led to crashes during planning if partitionwise joining is enabled.

- Disallow copying of invalidated replication slots (Shlok Kyal)

This prevents trouble when the invalid slot points to WAL that's already been removed.

- Disallow restoring logical replication slots on standby servers that are not in hot-standby mode (Masahiko Sawada)

This prevents a scenario where the slot could remain valid after promotion even if `wal_level` is too low.

- Prevent over-advancement of catalog xmin in “fast forward” mode of logical decoding (Zhijie Hou)

This mistake could allow deleted catalog entries to be vacuumed away even though they were still potentially needed by the WAL-reading process.

- Avoid data loss when DDL operations that don't take a strong lock affect tables that are being logically replicated (Shlok Kyal, Hayato Kuroda)

The catalog changes caused by the DDL command were not reflected into WAL-decoding processes, allowing them to decode subsequent changes using stale catalog data, probably resulting in data corruption.

- Prevent incorrect reset of replication origin when an apply worker encounters an error but the error is caught and does not result in worker exit (Hayato Kuroda)

This mistake could allow duplicate data to be applied.

- Fix crash in logical replication if the subscriber's partitioned table has a BRIN index (Tom Lane)
- Avoid duplicate snapshot creation in logical replication index lookups (Heikki Linnakangas)
- Improve detection of mixed-origin subscriptions (Hou Zhijie, Shlok Kyal)

Subscription creation gives a warning if a subscribed-to table is also being followed through other publications, since that could cause duplicate data to be received. This change improves that logic to also detect cases where a partition parent or child table is the one being followed through another publication.

- Fix wrong checkpoint details in error message about incorrect recovery timeline choice (David Steele)

If the requested recovery timeline is not reachable, the reported checkpoint and timeline should be the values read from the `backup_label`, if there is one. This message previously reported values

from the control file, which is correct when recovering from the control file without a `backup_label`, but not when there is a `backup_label`.

- Fix order of operations in `smgropen()` (Andres Freund)

Ensure that the `SMgrRelation` object is fully initialized before calling the `smgr_open` callback, so that it can be cleaned up properly if the callback fails.

- Remove incorrect assertion in `pgstat_report_stat()` (Michael Paquier)
- Fix overly-strict assertion in `gistFindCorrectParent()` (Heikki Linnakangas)
- Avoid assertion failure in parallel vacuum when `maintenance_work_mem` has a very small value (Masahiko Sawada)
- Fix rare assertion failure in standby servers when the primary is restarted (Heikki Linnakangas)
- In PL/pgSQL, avoid “unexpected plan node type” error when a scrollable cursor is defined on a simple `SELECT expression` query (Andrei Lepikhov)
- Don't try to drop individual index partitions in `pg_dump`'s `--clean` mode (Jian He)

The server rejects such `DROP` commands. That has no real consequences, since the partitions will go away anyway in the subsequent `DROPS` of either their parent tables or their partitioned index. However, the error reported for the attempted drop causes problems when restoring in `--single-transaction` mode.

- In `pg_dumpall`, avoid emitting invalid role `GRANT` commands if `pg_auth_members` contains invalid role OIDs (Tom Lane)

Instead, print a warning and skip the entry. This copes better with catalog corruption that has been seen to occur in back branches as a result of race conditions between `GRANT` and `DROP ROLE`.

- In `pg_amcheck` and `pg_upgrade`, use the correct function to free allocations made by `libpq` (Michael Paquier, Ranier Vilela)

These oversights could result in crashes in certain Windows build configurations, such as a debug build of `libpq` used by a non-debug build of the calling application.

- Fix `reindexdb`'s scheduling of parallel reindex operations (Alexander Korotkov)

The original coding failed to achieve the expected amount of parallelism.

- Avoid crashing with corrupt input data in `contrib/pageinspect`'s `heap_page_items()` (Dmitry Kovalenko)
- Prevent assertion failure in `contrib/pg_freespacemap`'s `pg_freespacemap()` (Tender Wang)

Applying `pg_freespacemap()` to a relation lacking storage (such as a view) caused an assertion failure, although there was no ill effect in non-assert builds. Add an error check to reject that case.

- In `contrib/postgres_fdw`, avoid pulling up restriction conditions from subqueries (Alexander Pyhalov)

This fix prevents rare cases of “unexpected expression in subquery output” errors.

- Fix build failure when an old version of `libpq-fe.h` is present in system include directories (Tom Lane)
- Fix build failure on macOS 15.4 (Tom Lane, Peter Eisentraut)

This macOS update broke our configuration probe for `strchrnul()`.

- Fix `valgrind` labeling of per-buffer data of read streams (Thomas Munro)

This affects no core code in released versions of PostgreSQL, but an extension using the per-buffer data feature might have encountered spurious failures when being tested under `valgrind`.

- Avoid valgrind complaints about string hashing code (John Naylor)
- Update time zone data files to tzdata release 2025b for DST law changes in Chile, plus historical corrections for Iran (Tom Lane)

There is a new time zone America/Coyhaique for Chile's Aysén Region, to account for it changing to UTC-03 year-round and thus diverging from America/Santiago.

E.3. Release 17.4

Release date: 2025-02-20

This release contains a few fixes from 17.3. For information about new features in major release 17, see [Section E.7](#).

E.3.1. Migration to Version 17.4

A dump/restore is not required for those running 17.X.

However, if you are upgrading from a version earlier than 17.1, see [Section E.6](#).

E.3.2. Changes

- Improve behavior of libpq's quoting functions (Andres Freund, Tom Lane)

The changes made for CVE-2025-1094 had one serious oversight: `PQescapeLiteral()` and `PQescapeIdentifier()` failed to honor their string length parameter, instead always reading to the input string's trailing null. This resulted in including unwanted text in the output, if the caller intended to truncate the string via the length parameter. With very bad luck it could cause a crash due to reading off the end of memory.

In addition, modify all these quoting functions so that when invalid encoding is detected, an invalid sequence is substituted for just the first byte of the presumed character, not all of it. This reduces the risk of problems if a calling application performs additional processing on the quoted string.

- Fix small memory leak in `pg_createsubscriber` (Ranier Vilela)
- Fix meson build system to correctly detect availability of the `bsd_auth.h` system header (Nazir Bilal Yavuz)

E.4. Release 17.3

Release date: 2025-02-13

This release contains a variety of fixes from 17.2. For information about new features in major release 17, see [Section E.7](#).

E.4.1. Migration to Version 17.3

A dump/restore is not required for those running 17.X.

However, if you are upgrading from a version earlier than 17.1, see [Section E.6](#).

E.4.2. Changes

- Harden `PQescapeString` and allied functions against invalidly-encoded input strings (Andres Freund, Noah Misch)

Data-quoting functions supplied by libpq now fully check the encoding validity of their input. If invalid characters are detected, they report an error if possible. For the ones that lack an error return convention, the output string is adjusted to ensure that the server will report invalid encod-

ing and no intervening processing will be fooled by bytes that might happen to match single quote, backslash, etc.

The purpose of this change is to guard against SQL-injection attacks that are possible if one of these functions is used to quote crafted input. There is no hazard when the resulting string is sent directly to a PostgreSQL server (which would check its encoding anyway), but there is a risk when it is passed through `psql` or other client-side code. Historically such code has not carefully vetted encoding, and in many cases it's not clear what it should do if it did detect such a problem.

This fix is effective only if the data-quoting function, the server, and any intermediate processing agree on the character encoding that's being used. Applications that insert untrusted input into SQL commands should take special care to ensure that that's true.

Applications and drivers that quote untrusted input without using these libpq functions may be at risk of similar problems. They should first confirm the data is valid in the encoding expected by the server.

The PostgreSQL Project thanks Stephen Fewer for reporting this problem. (CVE-2025-1094)

- Restore auto-truncation of database and user names appearing in connection requests (Nathan Bossart)

This reverts a v17 change that proved to cause trouble for some users. Over-length names should be truncated in an encoding-aware fashion, but for now just return to the former behavior of blind truncation at `NAMEDATALEN-1` bytes.

- Exclude parallel workers from connection privilege checks and limits (Tom Lane)

Do not check `datallowconn`, `rolcanlogin`, and `ACL_CONNECT` privileges when starting a parallel worker, instead assuming that it's enough for the leader process to have passed similar checks originally. This avoids, for example, unexpected failures of parallelized queries when the leader is running as a role that lacks login privilege. In the same vein, enforce `ReservedConnections`, `datconnlimit`, and `rolconnlimit` limits only against regular backends, and count only regular backends while checking if the limits were already reached. Those limits are meant to prevent excessive consumption of process slots for regular backends --- but parallel workers and other special processes have their own pools of process slots with their own limit checks.

- Drop “Lock” suffix from LWLock wait event names (Bertrand Drouvot)

Refactoring unintentionally caused the `pg_stat_activity` view to show lock-related wait event names with a “Lock” suffix, which among other things broke joining it to `pg_wait_events`.

- Fix possible failure to return all matching tuples for a btree index scan with a `ScalarArrayOp (= ANY)` condition (Peter Geoghegan)
- Fix possible re-use of stale results in window aggregates (David Rowley)

A window aggregate with a “run condition” optimization and a pass-by-reference result type might incorrectly return the result from the previous partition instead of performing a fresh calculation.

- Keep `TransactionXmin` in sync with `MyProc->xmin` (Heikki Linnakangas)

This oversight could permit a process to try to access data that had already been vacuumed away. One known consequence is transient “could not access status of transaction” errors.

- Fix race condition that could cause failure to add a newly-inserted catalog entry to a catalog cache list (Heikki Linnakangas)

This could result, for example, in failure to use a newly-created function within an existing session.

- Prevent possible catalog corruption when a system catalog is vacuumed concurrently with an update (Noah Misch)
- Fix data corruption when relation truncation fails (Thomas Munro)

The filesystem calls needed to perform relation truncation could fail, leaving inconsistent state on disk (for example, effectively reviving deleted data). We can't really prevent that, but we can recover by dint of making such failures into PANICs, so that consistency is restored by replaying from WAL up to just before the attempted truncation. This isn't a hugely desirable behavior, but such failures are rare enough that it seems an acceptable solution.

- Prevent checkpoints from starting during relation truncation (Robert Haas)

This avoids a race condition wherein the modified file might not get fsync'd before completing the checkpoint, creating a risk of data corruption if the operating system crashes soon after.

- Avoid possibly losing an update of `pg_database.datfrozenxid` when `VACUUM` runs concurrently with a `REASSIGN OWNED` that changes that database's owner (Kirill Reshke)
- Fix incorrect `tg_updatedcols` values passed to `AFTER UPDATE` triggers (Tom Lane)

In some cases the `tg_updatedcols` bitmap could describe the set of columns updated by an earlier command in the same transaction, fooling the trigger into doing the wrong thing.

Also, prevent memory bloat caused by making too many copies of the `tg_updatedcols` bitmap.

- Fix detach of a partition that has its own foreign-key constraint referencing a partitioned table (Amul Sul)

In common cases, foreign keys are defined on a partitioned table's top level; but if instead one is defined on a partition and references a partitioned table, and the referencing partition is detached, the relevant `pg_constraint` entries were updated incorrectly. This led to errors like “could not find ON INSERT check triggers of foreign key constraint”.

- Fix `pg_get_constraintdef`'s support for `NOT NULL` constraints on domains (Álvaro Herrera)
- Fix mis-processing of `to_timestamp`'s `FFn` format codes (Tom Lane)

An integer format code immediately preceding `FFn` would consume all available digits, leaving none for `FFn`.

- When deparsing a `PASSING` clause in a SQL/JSON query function, ensure that variable names are double-quoted when necessary (Dean Rasheed)
- When deparsing an `XMLTABLE()` expression, ensure that XML namespace names are double-quoted when necessary (Dean Rasheed)
- Include the `ldapscheme` option in `pg_hba_file_rules()` output (Laurenz Albe)
- Fix planning of pre-sorted `UNION` operations for cases where the input column datatypes don't all match (David Rowley)

This error could lead to sorting data with the wrong sort operator, with consequences ranging from no visible problem to core dumps.

- Don't merge `UNION` operations if their column collations aren't consistent (Tom Lane)

Previously we ignored collations when deciding if it's safe to merge `UNION` steps into a single N-way `UNION` operation. This was arguably valid before the introduction of nondeterministic collations, but it's not anymore, since the collation in use can affect the definition of uniqueness.

- Prevent “wrong varnullingrels” planner errors after pulling up a subquery that's underneath an outer join (Tom Lane)
- Ignore nulling-relation marker bits when looking up statistics (Richard Guo)

This oversight could lead to failure to use relevant statistics about expressions, or to “corrupt MVNDDistinct entry” errors.

- Fix missed expression processing for partition pruning steps (Tom Lane)

This oversight could lead to “unrecognized node type” errors, and perhaps other problems, in queries accessing partitioned tables.

- Give the slotsync worker process its own process slot (Tom Lane, Hou Zhijie)

This was overlooked in the addition of the slotsync worker, with the result that its process slot effectively came out of the pool meant for regular backend processes. This could result in failure to launch the worker, or to subsequent failures of connection requests that should have succeeded according to the configured settings, if the number of regular backend processes approached `max_connections`.

- Allow dshash tables to grow past 1GB (Matthias van de Meent)

This avoids errors like “invalid DSA memory alloc request size”. The case can occur for example in transactions that process several million tables.

- Avoid possible integer overflow in `bringetbitmap()` (James Hunter, Evgeniy Gorbanyov)

Since the result is only used for statistical purposes, the effects of this error were mostly cosmetic.

- Correct miscalculation of SLRU bank numbers (Yura Sokolov)

This error led to using a smaller number of banks than intended, causing more contention but no functional misbehavior.

- Ensure that an already-set process latch doesn't prevent the postmaster from noticing socket events (Thomas Munro)

An extremely heavy workload of backends launching workers and workers exiting could prevent the postmaster from responding to incoming client connections in a timely fashion.

- Prevent streaming standby servers from looping infinitely when reading a WAL record that crosses pages (Kyotaro Horiguchi, Alexander Kukushkin)

This would happen when the record's continuation is on a page that needs to be read from a different WAL source.

- Fix unintended promotion of FATAL errors to PANIC during early process startup (Noah Misch)

This fixes some unlikely cases that would result in “PANIC: `proc_exit()` called in child process”.

- Fix cases where an operator family member operator or support procedure could become a dangling reference (Tom Lane)

In some cases a data type could be dropped while references to its OID still remain in `pg_amop` or `pg_amproc`. While that caused no immediate issues, an attempt to drop the owning operator family would fail, and `pg_dump` would produce bogus output when dumping the operator family. This fix causes creation and modification of operator families/classes to add needed dependency entries so that dropping a data type will also drop any dependent operator family elements. That does not help vulnerable pre-existing operator families, though, so a band-aid has also been added to `DROP OPERATOR FAMILY` to prevent failure when dropping a family that has dangling members.

- Fix multiple memory leaks in logical decoding output (Vignesh C, Masahiko Sawada, Boyu Yang)
- Fix small memory leak when updating the `application_name` or `cluster_name` settings (Tofig Aliev)
- Avoid crash when a background process tries to check a new value of `synchronized_standby_slots` (Álvaro Herrera)
- Avoid integer overflow while testing `wal_skip_threshold` condition (Tom Lane)

A transaction that created a very large relation could mistakenly decide to ensure durability by copying the relation into WAL instead of fsync'ing it, thereby negating the point of

`wal_skip_threshold`. (This only matters when `wal_level` is set to `minimal`, else a WAL copy is required anyway.)

- Fix unsafe order of operations during cache lookups (Noah Misch)

The only known consequence was a usually-harmless “you don't own a lock of type `ExclusiveLock`” warning during `GRANT TABLESPACE`.

- Avoid potential use-after-free in parallel vacuum (Vallimaharajan G, John Naylor)

This bug seems to have no consequences in standard builds, but it's theoretically a hazard.

- Fix possible “failed to resolve name” failures when using JIT on older ARM platforms (Thomas Munro)

This could occur as a consequence of inconsistency about the default setting of `-moutline-atomics` between `gcc` and `clang`. At least Debian and Ubuntu are known to ship `gcc` and `clang` compilers that target `armv8-a` but differ on the use of outline atomics by default.

- Fix assertion failure in `WITH RECURSIVE ... UNION` queries (David Rowley)
- Avoid assertion failure in rule deparsing if a set operation leaf query contains set operations (Man Zeng, Tom Lane)
- Avoid edge-case assertion failure in parallel query startup (Tom Lane)
- Fix assertion failure at shutdown when writing out the statistics file (Michael Paquier)
- Avoid `valgrind` complaints about string hashing code (John Naylor)
- In `NULLIF()`, avoid passing a read-write expanded object pointer to the data type's equality function (Tom Lane)

The equality function could modify or delete the object if it's given a read-write pointer, which would be bad if we decide to return it as the `NULLIF()` result. There is probably no problem with any built-in equality function, but it's easy to demonstrate a failure with one coded in PL/pgSQL.

- Ensure that expression preprocessing is applied to a default null value in `INSERT` (Tom Lane)

If the target column is of a domain type, the planner must insert a `coerce-to-domain` step not just a null constant, and this expression missed going through some required processing steps. There is no known consequence with domains based on core data types, but in theory an error could occur with domains based on extension types.

- Avoid data loss when starting a bulk write on a relation fork that already contains data (Matthias van de Meent)

Any pre-existing data was overwritten with zeroes. This is not an issue for core PostgreSQL, which never does that. Some extensions would like to, however.

- Avoid crash if a server process tried to iterate over a shared radix tree that it didn't create (Masahiko Sawada)

There is no code in core PostgreSQL that does this, but an extension might wish to.

- Repair memory leaks in PL/Python (Mat Arye, Tom Lane)

Repeated use of `PLyPlan.execute` or `plpy.cursor` resulted in memory leakage for the duration of the calling PL/Python function.

- Fix PL/Tcl to compile with Tcl 9 (Peter Eisentraut)
- In the `ecpg` preprocessor, fix possible misprocessing of cursors that reference out-of-scope variables (Tom Lane)
- In `ecpg`, fix compile-time warnings about unsupported use of `COPY ... FROM STDIN` (Ryo Kanbayashi)

Previously, the intended warning was not issued due to a typo.

- Fix `psql` to safely handle file path names that are encoded in SJIS (Tom Lane)

Some two-byte characters in SJIS have a second byte that is equal to ASCII backslash (`\`). These characters were corrupted by path name normalization, preventing access to files whose names include such characters.

- Add `psql` tab completion for `COPY (MERGE INTO)` (Jian He)
- Fix use of wrong version of `pqsignal()` in `pgbench` and `psql` (Fujii Masao, Tom Lane)

This error could lead to misbehavior when using the `-T` option in `pgbench` or the `\watch` command in `psql`, due to interrupted system calls not being resumed as expected.

- Fix misexecution of some nested `\if` constructs in `pgbench` (Michail Nikolaev)

An `\if` command appearing within a false (not-being-executed) `\if` branch was incorrectly treated the same as `\elif`.

- In `pgbench`, fix possible misdisplay of progress messages during table initialization (Yushi Ogiwara, Tatsuo Ishii, Fujii Masao)
- Make `pg_controldata` more robust against corrupted `pg_control` files (Ilyasov Ian, Anton Voloshin)

Since `pg_controldata` will attempt to print the contents of `pg_control` even if the CRC check fails, it must take care not to misbehave for invalid field values. This patch fixes some issues triggered by invalid timestamps and apparently-negative WAL segment sizes.

- Fix possible crash in `pg_dump` with identity sequences attached to tables that are extension members (Tom Lane)
- Fix memory leak in `pg_restore` with `zstd`-compressed data (Tom Lane)

The leak was per-decompression-operation, so would be most noticeable with a dump containing many tables or large objects.

- Fix `pg_basebackup` to correctly handle `pg_wal.tar` files exceeding 2GB on Windows (Davinder Singh, Thomas Munro)
- Use SQL-standard function bodies in the declarations of `contrib/earthdistance`'s SQL-language functions (Tom Lane, Ronan Dunklau)

This change allows their references to `contrib/cube` to be resolved during extension creation, reducing the risk of search-path-based failures and possible attacks.

In particular, this restores their usability in contexts like generated columns, for which PostgreSQL v17 restricts the search path on security grounds. We have received reports of databases failing to be upgraded to v17 because of that. This patch has been included in v16 to provide a workaround: updating the `earthdistance` extension to this version beforehand should allow an upgrade to succeed.

- Detect version mismatch between `contrib/pageinspect`'s SQL declarations and the underlying shared library (Tomas Vondra)

Previously, such a mismatch could result in a crash while calling `brin_page_items()`. Instead throw an error recommending updating the extension.

- When trying to cancel a remote query in `contrib/postgres_fdw`, re-issue the cancel request a few times if it didn't seem to do anything (Tom Lane)

This fixes a race condition where we might try to cancel a just-sent query before the remote server has started to process it, so that the initial cancel request is ignored.

- Update configuration probes that determine the compiler switches needed to access ARM CRC instructions (Tom Lane)

On ARM platforms where the baseline CPU target lacks CRC instructions, we need to supply a `-march` switch to persuade the compiler to compile such instructions. Recent versions of `gcc` reject the value we were trying, leading to silently falling back to software CRC.

- Fix meson build system to support old OpenSSL libraries on Windows (Darek Slusarczyk)

Add support for the legacy library names `ssleay32` and `libeay32`.

- In Windows builds using meson, ensure all `libcommon` and `libpgport` functions are exported (Vladlen Popolitov, Heikki Linnakangas)

This fixes “unresolved external symbol” build errors for extensions.

- Fix meson configuration process to correctly detect OSSP's `uuid.h` header file under MSVC (Andrew Dunstan)
- When building with meson, install `pgevent` in `pkglibdir` not `bindir` (Peter Eisentraut)

This matches the behavior of the make-based build system and the old MSVC build system.

- When building with meson, install `sepgsql.sql` under `share/contrib/` not `share/extension/` (Peter Eisentraut)

This matches what the make-based build system does.

- Update time zone data files to tzdata release 2025a for DST law changes in Paraguay, plus historical corrections for the Philippines (Tom Lane)

E.5. Release 17.2

Release date: 2024-11-21

This release contains a few fixes from 17.1. For information about new features in major release 17, see [Section E.7](#).

E.5.1. Migration to Version 17.2

A dump/restore is not required for those running 17.X.

However, if you are upgrading from a version earlier than 17.1, see [Section E.6](#).

E.5.2. Changes

- Repair ABI break for extensions that work with struct `ResultRelInfo` (Tom Lane)

Last week's minor releases unintentionally broke binary compatibility with `timescaledb` and several other extensions. Restore the affected structure to its previous size, so that such extensions need not be rebuilt.

- Restore functionality of `ALTER {ROLE|DATABASE} SET role` (Tom Lane, Noah Misch)

The fix for CVE-2024-10978 accidentally caused settings for `role` to not be applied if they come from non-interactive sources, including previous `ALTER {ROLE|DATABASE}` commands and the `PGOPTIONS` environment variable.

- Fix cases where a logical replication slot's `restart_lsn` could go backwards (Masahiko Sawada)

Previously, restarting logical replication could sometimes cause the slot's restart point to be recomputed as an older value than had previously been advertised in `pg_replication_slots`. This is bad, since for example WAL files might have been removed on the basis of the later `restart_lsn` value, in which case replication would fail to restart.

- Avoid deleting still-needed WAL files during `pg_rewind` (Polina Bungina, Alexander Kukushkin)

Previously, in unlucky cases, it was possible for `pg_rewind` to remove important WAL files from the rewound demoted primary. In particular this happens if those files have been marked for archival (i.e., their `.ready` files were created) but not yet archived. Then the newly promoted node no longer has such files because of them having been recycled, but likely they are needed for recovery in the demoted node. If `pg_rewind` removes them, recovery is not possible anymore.

- Fix race conditions associated with dropping shared statistics entries (Kyotaro Horiguchi, Michael Paquier)

These bugs could lead to loss of statistics data, assertion failures, or “can only drop stats once” errors.

- Count index scans in `contrib/bloom` indexes in the statistics views, such as the `pg_stat_user_indexes.idx_scan` counter (Masahiro Ikeda)
- Fix crash when checking to see if an index's opclass options have changed (Alexander Korotkov)

Some forms of `ALTER TABLE` would fail if the table has an index with non-default operator class options.

- Avoid assertion failure caused by disconnected NFA sub-graphs in regular expression parsing (Tom Lane)

This bug does not appear to have any visible consequences in non-assert builds.

E.6. Release 17.1

Release date: 2024-11-14

This release contains a variety of fixes from 17.0. For information about new features in major release 17, see [Section E.7](#).

E.6.1. Migration to Version 17.1

A dump/restore is not required for those running 17.X.

However, if you have ever detached a partition from a partitioned table that has a foreign-key reference to another partitioned table, and not dropped the former partition, then you may have catalog and/or data corruption to repair, as detailed in the fifth changelog entry below.

Also, in the uncommon case that a database's `LC_CTYPE` setting is `C` while its `LC_COLLATE` setting is some other locale, indexes on textual columns should be reindexed, as described in the sixth changelog entry below.

E.6.2. Changes

- Ensure cached plans are marked as dependent on the calling role when RLS applies to a non-top-level table reference (Nathan Bossart)

If a CTE, subquery, sublink, security invoker view, or coercion projection in a query references a table with row-level security policies, we neglected to mark the resulting plan as potentially dependent on which role is executing it. This could lead to later query executions in the same session using the wrong plan, and then returning or hiding rows that should have been hidden or returned instead.

The PostgreSQL Project thanks Wolfgang Walther for reporting this problem. (CVE-2024-10976)

- Make libpq discard error messages received during SSL or GSS protocol negotiation (Jacob Champion)

An error message received before encryption negotiation is completed might have been injected by a man-in-the-middle, rather than being real server output. Reporting it opens the door to var-

ious security hazards; for example, the message might spoof a query result that a careless user could mistake for correct output. The best answer seems to be to discard such data and rely only on libpq's own report of the connection failure.

The PostgreSQL Project thanks Jacob Champion for reporting this problem. (CVE-2024-10977)

- Fix unintended interactions between `SET SESSION AUTHORIZATION` and `SET ROLE` (Tom Lane)

The SQL standard mandates that `SET SESSION AUTHORIZATION` have a side-effect of doing `SET ROLE NONE`. Our implementation of that was flawed, creating more interaction between the two settings than intended. Notably, rolling back a transaction that had done `SET SESSION AUTHORIZATION` would revert `ROLE` to `NONE` even if that had not been the previous state, so that the effective user ID might now be different from what it had been before the transaction. Transiently setting `session_authorization` in a function `SET` clause had a similar effect. A related bug was that if a parallel worker inspected `current_setting('role')`, it saw `none` even when it should see something else.

The PostgreSQL Project thanks Tom Lane for reporting this problem. (CVE-2024-10978)

- Prevent trusted PL/Perl code from changing environment variables (Andrew Dunstan, Noah Misch)

The ability to manipulate process environment variables such as `PATH` gives an attacker opportunities to execute arbitrary code. Therefore, “trusted” PLs must not offer the ability to do that. To fix `plperl`, replace `%ENV` with a tied hash that rejects any modification attempt with a warning. Untrusted `plperl` retains the ability to change the environment.

The PostgreSQL Project thanks Coby Abrams for reporting this problem. (CVE-2024-10979)

- Fix updates of catalog state for foreign-key constraints when attaching or detaching table partitions (Jehan-Guillaume de Rorthais, Tender Wang, Álvaro Herrera)

If the referenced table is partitioned, then different catalog entries are needed for a referencing table that is stand-alone versus one that is a partition. `ATTACH/DETACH PARTITION` commands failed to perform this conversion correctly. In particular, after `DETACH` the now stand-alone table would be missing foreign-key enforcement triggers, which could result in the table later containing rows that fail the foreign-key constraint. A subsequent `re-ATTACH` could fail with surprising errors, too.

The way to fix this is to do `ALTER TABLE DROP CONSTRAINT` on the now stand-alone table for each faulty constraint, and then re-add the constraint. If re-adding the constraint fails, then some erroneous data has crept in. You will need to manually re-establish consistency between the referencing and referenced tables, then re-add the constraint.

This query can be used to identify broken constraints and construct the commands needed to recreate them:

```
SELECT conrelid::pg_catalog.regclass AS "constrained table",
       conname AS constraint,
       confrelid::pg_catalog.regclass AS "references",
       pg_catalog.format('ALTER TABLE %s DROP CONSTRAINT %I;',
                          conrelid::pg_catalog.regclass, conname) AS "drop",
       pg_catalog.format('ALTER TABLE %s ADD CONSTRAINT %I %s;',
                          conrelid::pg_catalog.regclass, conname,
                          pg_catalog.pg_get_constraintdef(oid)) AS "add"
FROM pg_catalog.pg_constraint c
WHERE contype = 'f' AND conparentid = 0 AND
      (SELECT count(*) FROM pg_catalog.pg_constraint c2
       WHERE c2.conparentid = c.oid) <>
      (SELECT count(*) FROM pg_catalog.pg_inherits i
       WHERE (i.inhparent = c.conrelid OR i.inhparent = c.confrelid) AND
            EXISTS (SELECT 1 FROM pg_catalog.pg_partitioned_table
                    WHERE partrelid = i.inhparent));
```

Since it is possible that one or more of the `ADD CONSTRAINT` steps will fail, you should save the query's output in a file and then attempt to perform each step.

- Fix test for C locale when `LC_COLLATE` is different from `LC_CTYPE` (Jeff Davis)

When using `libc` as the default collation provider, the test to see if C locale is in use for collation accidentally checked `LC_CTYPE` not `LC_COLLATE`. This has no impact in the typical case where those settings are the same, nor if both are not C (nor its alias `POSIX`). However, if `LC_CTYPE` is C while `LC_COLLATE` is some other locale, wrong query answers could ensue, and corruption of indexes on strings was possible. Users of databases with such settings should reindex affected indexes after installing this update. The converse case with `LC_COLLATE` being C while `LC_CTYPE` is some other locale would cause performance degradation, but no actual errors.

- Don't use partitionwise joins or grouping if the query's collation for the key column doesn't match the partition key's collation (Jian He, Webbo Han)

Such plans could produce incorrect results.

- Avoid planner failure after converting an `IS NULL` test on a `NOT NULL` column to constant `FALSE` (Richard Guo)

This bug typically led to errors such as “variable not found in subplan target lists”.

- Avoid possible planner crash while inlining a SQL function whose arguments contain certain array-related constructs (Tom Lane, Nathan Bossart)
- Fix possible wrong answers or “wrong varnullingrels” planner errors for `MERGE ... WHEN NOT MATCHED BY SOURCE` actions (Dean Rasheed)
- Fix possible “could not find pathkey item to sort” error when the output of a `UNION ALL` member query needs to be sorted, and the sort column is an expression (Andrei Lepikhov, Tom Lane)
- Fix edge case in B-tree `ScalarArrayOp` index scans (Peter Geoghegan)

When a scrollable cursor with a plan of this kind was backed up to its starting point and then run forward again, wrong answers were possible.

- Fix assertion failure or confusing error message for `COPY (query) TO ...`, when the *query* is rewritten by a `DO INSTEAD NOTIFY` rule (Tender Wang, Tom Lane)
- Fix validation of `COPY`'s `FORCE_NOT_NULL` and `FORCE_NULL` options (Joel Jacobson)

Some incorrect usages are now rejected as they should be.

- Fix server crash when a `json_objectagg()` call contains a volatile function (Amit Langote)
- Fix detection of skewed data during parallel hash join (Thomas Munro)

After repartitioning the inner side of a hash join because one partition has accumulated too many tuples, we check to see if all the partition's tuples went into the same child partition, which suggests that they all have the same hash value and further repartitioning cannot improve matters. This check malfunctioned in some cases, allowing repeated futile repartitioning which would eventually end in a resource-exhaustion error.

- Avoid crash when `ALTER DATABASE SET` is used to set a server parameter that requires search-path-based lookup, such as `default_text_search_config` (Jeff Davis)
- Avoid repeated lookups of opclasses and collations while creating a new index on a partitioned table (Tom Lane)

This was problematic mainly because some of the lookups would be done with a restricted `search_path`, leading to unexpected failures if the `CREATE INDEX` command referenced objects outside `pg_catalog`.

This fix also prevents comments on the parent partitioned index from being copied to child indexes.

- Add missing dependency from a partitioned table to a non-built-in access method specified in `CREATE TABLE ... USING` (Michael Paquier)

Dropping the access method should be blocked when a table exists that depends on it, but it was not, allowing subsequent odd behavior. Note that this fix only prevents problems for partitioned tables created after this update.

- Disallow locale names containing non-ASCII characters (Thomas Munro)

This is only an issue on Windows, as such locale names are not used elsewhere. They are problematic because it's quite unclear what encoding such names are represented in (since the locale itself defines the encoding to use). In recent PostgreSQL releases, an abort in the Windows runtime library could occur because of confusion about that.

Anyone who encounters the new error message should either create a new duplicated locale with an ASCII-only name using Windows Locale Builder, or consider using BCP 47-compliant locale names like `tr-TR`.

- Fix race condition in committing a serializable transaction (Heikki Linnakangas)

Mis-processing of a recently committed transaction could lead to an assertion failure or a “could not access status of transaction” error.

- Fix race condition in `COMMIT PREPARED` that resulted in orphaned 2PC files (wuchengwen)

A concurrent `PREPARE TRANSACTION` could cause `COMMIT PREPARED` to not remove the on-disk two-phase state file for the completed transaction. There was no immediate ill effect, but a subsequent crash-and-recovery could fail with “could not access status of transaction”, requiring manual removal of the orphaned file to restore service.

- Avoid invalid memory accesses after skipping an invalid toast index during `VACUUM FULL` (Tender Wang)

A list tracking yet-to-be-rebuilt indexes was not properly updated in this code path, risking assertion failures or crashes later on.

- Fix ways in which an “in place” catalog update could be lost (Noah Misch)

Normal row updates write a new version of the row to preserve rollback-ability of the transaction. However, certain system catalog updates are intentionally non-transactional and are done with an in-place update of the row. These patches fix race conditions that could cause the effects of an in-place update to be lost. As an example, it was possible to forget having set `pg_class.relhasindex` to true, preventing updates of the new index and thus causing index corruption.

- Reset catalog caches at end of recovery (Noah Misch)

This prevents scenarios wherein an in-place catalog update could be lost due to using stale data from a catalog cache.

- Avoid using parallel query while holding off interrupts (Francesco Degrossi, Noah Misch, Tom Lane)

This situation cannot arise normally, but it can be reached with test scenarios such as using a SQL-language function as B-tree support (which would be far too slow for production usage). If it did occur it would result in an indefinite wait.

- Ignore not-yet-defined Portals in the `pg_cursors` view (Tom Lane)

It is possible for user-defined code that inspects this view to be called while a new cursor is being set up, and if that happens a null pointer dereference would ensue. Avoid the problem by defining the view to exclude incompletely-set-up cursors.

- Avoid “unexpected table_index_fetch_tuple call during logical decoding” error while decoding a transaction involving insertion of a column default value (Takeshi Ideriha, Hou Zhijie)

- Reduce memory consumption of logical decoding (Masahiko Sawada)

Use a smaller default block size to store tuple data received during logical replication. This reduces memory wastage, which has been reported to be severe while processing long-running transactions, even leading to out-of-memory failures.

- Fix behavior of stable functions called from a `CALL` statement's argument list, when the `CALL` is within a PL/pgSQL `EXCEPTION` block (Tom Lane)

As with a similar fix in our previous quarterly releases, this case allowed such functions to be passed the wrong snapshot, causing them to see stale values of rows modified since the start of the outer transaction.

- Parse libpq's `keepalives` connection option in the same way as other integer-valued options (Yuto Sasaki)

The coding used here rejected trailing whitespace in the option value, unlike other cases. This turns out to be problematic in `ecpg`'s usage, for example.

- In `ecpglib`, fix out-of-bounds read when parsing incorrect datetime input (Bruce Momjian, Pavel Nekrasov)

It was possible to try to read the location just before the start of a constant array. Real-world consequences seem minimal, though.

- Fix `psql`'s `describe` commands to again work with pre-9.4 servers (Tom Lane)

Commands involving display of an ACL (permissions) column failed with very old PostgreSQL servers, due to use of a function not present in those versions.

- Avoid hanging if an interval less than 1ms is specified in `psql`'s `\watch` command (Andrey Borodin, Michael Paquier)

Instead, treat this the same as an interval of zero (no wait between executions).

- Fix failure to find replication password in `~/.pgpass` (Tom Lane)

`pg_basebackup` and `pg_receivewal` failed to match an entry in `~/.pgpass` that had `replication` in the database name field, if no `-d` or `--dbname` switch was supplied. This resulted in an unexpected prompt for password.

- In `pg_combinebackup`, throw an error if an incremental backup file is present in a directory that is supposed to contain a full backup (Robert Haas)
- In `pg_combinebackup`, don't construct filenames containing double slashes (Robert Haas)

This caused no functional problems, but the duplicate slashes were visible in error messages, which could create confusion.

- Avoid trying to reindex temporary tables and indexes in `vacuumdb` and in `parallel reindexdb` (VaibhaveS, Michael Paquier, Fujii Masao, Nathan Bossart)

Reindexing other sessions' temporary tables cannot work, but the check to skip them was missing in some code paths, leading to unwanted failures.

- Fix incorrect LLVM-generated code on ARM64 platforms (Thomas Munro, Anthonin Bonnefoy)

When using JIT compilation on ARM platforms, the generated code could not support relocation distances exceeding 32 bits, allowing unlucky placement of generated code to cause server crashes on large-memory systems.

- Fix a few places that assumed that process start time (represented as a `time_t`) will fit into a `long` value (Max Johnson, Nathan Bossart)

On platforms where `long` is 32 bits (notably Windows), this coding would fail after Y2038. Most of the failures appear only cosmetic, but notably `pg_ctl start` would hang.

- Update time zone data files to tzdata release 2024b (Tom Lane)

This tzdata release changes the old System-V-compatibility zone names to duplicate the corresponding geographic zones; for example `PST8PDT` is now an alias for `America/Los_Angeles`. The main visible consequence is that for timestamps before the introduction of standardized time zones, the zone is considered to represent local mean solar time for the named location. For example, in `PST8PDT`, `timestampz` input such as `1801-01-01 00:00` would previously have been rendered as `1801-01-01 00:00:00-08`, but now it is rendered as `1801-01-01 00:00:00-07:52:58`.

Also, historical corrections for Mexico, Mongolia, and Portugal. Notably, `Asia/Choibalsan` is now an alias for `Asia/Ulaanbaatar` rather than being a separate zone, mainly because the differences between those zones were found to be based on untrustworthy data.

E.7. Release 17

Release date: 2024-09-26

E.7.1. Overview

PostgreSQL 17 contains many new features and enhancements, including:

- New memory management system for `VACUUM`, which reduces memory consumption and can improve overall vacuuming performance.
- New SQL/JSON capabilities, including constructors, identity functions, and the `JSON_TABLE()` function, which converts JSON data into a table representation.
- Various query performance improvements, including for sequential reads using streaming I/O, write throughput under high concurrency, and searches over multiple values in a `btree` index.
- Logical replication enhancements, including:
 - Failover control
 - `pg_createsubscriber`, a utility that creates logical replicas from physical standbys
 - `pg_upgrade` now preserves logical replication slots on publishers and full subscription state on subscribers. This will allow upgrades to future major versions to continue logical replication without requiring copy to resynchronize.
- New client-side connection option, `sslnegotiation=direct`, that performs a direct TLS handshake to avoid a round-trip negotiation.
- `pg_basebackup` now supports incremental backup.
- `COPY` adds a new option, `ON_ERROR ignore`, that allows a copy operation to continue in the event of an error.

The above items and other new features of PostgreSQL 17 are explained in more detail in the sections below.

E.7.2. Migration to Version 17

A dump/restore using `pg_dumpall` or use of `pg_upgrade` or logical replication is required for those wishing to migrate data from any previous release. See [Section 18.6](#) for general information on migrating to new major releases.

Version 17 contains a number of changes that may affect compatibility with previous releases. Observe the following incompatibilities:

- Change functions to use a safe `search_path` during maintenance operations (Jeff Davis)

This prevents maintenance operations (`ANALYZE`, `CLUSTER`, `CREATE INDEX`, `CREATE MATERIALIZED VIEW`, `REFRESH MATERIALIZED VIEW`, `REINDEX`, or `VACUUM`) from performing unsafe access. Functions

used by expression indexes and materialized views that need to reference non-default schemas must specify a search path during function creation.

- Restrict `ago` to only appear at the end in `interval` values (Joseph Koshakow)

Also, prevent empty interval units from appearing multiple times.

- Remove server variable `old_snapshot_threshold` (Thomas Munro)

This variable allowed vacuum to remove rows that potentially could be still visible to running transactions, causing "snapshot too old" errors later if accessed. This feature might be re-added to PostgreSQL later if an improved implementation is found.

- Change `SET SESSION AUTHORIZATION` handling of the initial session user's superuser status (Joseph Koshakow)

The new behavior is based on the session user's superuser status at the time the `SET SESSION AUTHORIZATION` command is issued, rather than their superuser status at connection time.

- Remove feature which simulated per-database users (Nathan Bossart)

The feature, `db_user_namespace`, was rarely used.

- Remove `adminpack` contrib extension (Daniel Gustafsson)

This was used by now end-of-life pgAdmin III.

- Remove `wal_sync_method` value `fsync_writethrough` on Windows (Thomas Munro)

This value was the same as `fsync` on Windows.

- Change file boundary handling of two WAL file name functions (Kyotaro Horiguchi, Andres Freund, Bruce Momjian)

The functions `pg_walfile_name()` and `pg_walfile_name_offset()` used to report the previous LSN segment number when the LSN was on a file segment boundary; it now returns the current LSN segment.

- Remove server variable `trace_recovery_messages` since it is no longer needed (Bharath Rupireddy)
- Remove `information schema` column `element_types.domain_default` (Peter Eisentraut)
- Change `pgrowlocks` lock mode output labels (Bruce Momjian)
- Remove `buffers_backend` and `buffers_backend_fsync` from `pg_stat_bgwriter` (Bharath Rupireddy)

These fields are considered redundant to similar columns in `pg_stat_io`.

- Rename I/O block read/write timing statistics columns of `pg_stat_statements` (Nazir Bilal Yavuz)

This renames `blk_read_time` to `shared_blk_read_time`, and `blk_write_time` to `shared_blk_write_time`.

- Change `pg_attribute.attstattarget` and `pg_statistic_ext.stxstattarget` to represent the default statistics target as `NULL` (Peter Eisentraut)
- Rename `pg_collation.colliculocale` to `colllocale` and `pg_database.daticulocale` to `datlocale` (Jeff Davis)
- Rename `pg_stat_progress_vacuum` column `max_dead_tuples` to `max_dead_tuple_bytes`, rename `num_dead_tuples` to `num_dead_item_ids`, and add `dead_tuple_bytes` (Masahiko Sawada)
- Rename SLRU columns in system view `pg_stat_slru` (Alvaro Herrera)

The column names accepted by `pg_stat_reset_slru()` are also changed.

E.7.3. Changes

Below you will find a detailed account of the changes between PostgreSQL 17 and the previous major release.

E.7.3.1. Server

E.7.3.1.1. Optimizer

- Allow the optimizer to improve CTE plans by considering the statistics and sort order of columns referenced in earlier row output clauses (Jian Guo, Richard Guo, Tom Lane)
- Improve optimization of `IS NOT NULL` and `IS NULL` query restrictions (David Rowley, Richard Guo, Andy Fan)

Remove `IS NOT NULL` restrictions from queries on `NOT NULL` columns and eliminate scans on `NOT NULL` columns if `IS NULL` is specified.

- Allow partition pruning on boolean columns on `IS [NOT] UNKNOWN` conditionals (David Rowley)
- Improve optimization of range values when using containment operators `<@` and `@>` (Kim Johan Andersson, Jian He)
- Allow correlated `IN` subqueries to be transformed into joins (Andy Fan, Tom Lane)
- Improve optimization of the `LIMIT` clause on partitioned tables, inheritance parents, and `UNION ALL` queries (Andy Fan, David Rowley)
- Allow query nodes to be run in parallel in more cases (Tom Lane)
- Allow `GROUP BY` columns to be internally ordered to match `ORDER BY` (Andrei Lepikhov, Teodor Sigaev)

This can be disabled using server variable [enable_group_by_reordering](#).

- Allow `UNION` (without `ALL`) to use MergeAppend (David Rowley)
- Fix MergeAppend plans to more accurately compute the number of rows that need to be sorted (Alexander Kuzmenkov)
- Allow [GiST](#) and [SP-GiST](#) indexes to be part of incremental sorts (Miroslav Bendik)

This is particularly useful for `ORDER BY` clauses where the first column has a GiST and SP-GiST index, and other columns do not.

- Add columns to `pg_stats` to report range-type histogram information (Egor Rogov, Soumyadeep Chakraborty)

E.7.3.1.2. Indexes

- Allow [btree](#) indexes to more efficiently find a set of values, such as those supplied by `IN` clauses using constants (Peter Geoghegan, Matthias van de Meent)
- Allow [BRIN](#) indexes to be created using parallel workers (Tomas Vondra, Matthias van de Meent)

E.7.3.1.3. General Performance

- Allow vacuum to more efficiently remove and freeze tuples (Melanie Plageman, Heikki Linnakangas)

WAL traffic caused by vacuum is also more compact.

- Allow vacuum to more efficiently store tuple references (Masahiko Sawada, John Naylor)

Additionally, vacuum is no longer silently limited to one gigabyte of memory when [maintenance_work_mem](#) or [autovacuum_work_mem](#) are higher.

- Optimize vacuuming of relations with no indexes (Melanie Plageman)

- Increase default `vacuum_buffer_usage_limit` to 2MB (Thomas Munro)
- Improve performance when checking roles with many memberships (Nathan Bossart)
- Improve performance of heavily-contended WAL writes (Bharath Rupireddy)
- Improve performance when transferring large blocks of data to a client (Melih Mutlu)
- Allow the grouping of file system reads with the new system variable `io_combine_limit` (Thomas Munro, Andres Freund, Melanie Plageman, Nazir Bilal Yavuz)

E.7.3.1.4. Monitoring

- Create system view `pg_stat_checkpoint` (Bharath Rupireddy, Anton A. Melnikov, Alexander Korotkov)

Relevant columns have been removed from `pg_stat_bgwriter` and added to this new system view.

- Improve control over resetting statistics (Atsushi Torikoshi, Bharath Rupireddy)

Allow `pg_stat_reset_shared()` (with no arguments) and `pg_stat_reset_shared(NULL)` to reset all shared statistics. Allow `pg_stat_reset_shared('slru')` and `pg_stat_reset_slru()` (with no arguments) to reset SLRU statistics, which was already possible with `pg_stat_reset_slru(NULL)`.

- Add log messages related to WAL recovery from backups (Andres Freund)
- Add `log_connections` log line for `trust` connections (Jacob Champion)
- Add log message to report walsender acquisition and release of replication slots (Bharath Rupireddy)

This is enabled by the server variable `log_replication_commands`.

- Add system view `pg_wait_events` that reports wait event types (Bertrand Drouvot)

This is useful for adding descriptions to wait events reported in `pg_stat_activity`.

- Add `wait events` for checkpoint delays (Thomas Munro)
- Allow vacuum to report the progress of index processing (Sami Imseih)

This appears in system view `pg_stat_progress_vacuum` columns `indexes_total` and `indexes_processed`.

E.7.3.1.5. Privileges

- Allow granting the right to perform maintenance operations (Nathan Bossart)

The permission can be granted on a per-table basis using the `MAINTAIN` privilege and on a per-role basis via the `pg_maintain` predefined role. Permitted operations are `VACUUM`, `ANALYZE`, `REINDEX`, `REFRESH MATERIALIZED VIEW`, `CLUSTER`, and `LOCK TABLE`.

- Allow roles with `pg_monitor` membership to execute `pg_current_logfile()` (Pavlo Golub, Nathan Bossart)

E.7.3.1.6. Server Configuration

- Add system variable `allow_alter_system` to disallow `ALTER SYSTEM` (Jelte Fennema-Nio, Gabriele Bartolini)
- Allow `ALTER SYSTEM` to set unrecognized custom server variables (Tom Lane)

This is also possible with `GRANT ON PARAMETER`.

- Add server variable `transaction_timeout` to restrict the duration of transactions (Andrey Borodin, Japin Li, Junwang Zhao, Alexander Korotkov)
- Add a builtin platform-independent collation provider (Jeff Davis)

This supports `C` and `C.UTF-8` collations.

- Add server variable `huge_pages_status` to report the use of huge pages by Postgres (Justin Pryzby)

This is useful when `huge_pages` is set to `try`.

- Add server variable to disable event triggers (Daniel Gustafsson)

The setting, `event_triggers`, allows for the temporary disabling of event triggers for debugging.

- Allow the `SLRU` cache sizes to be configured (Andrey Borodin, Dilip Kumar, Alvaro Herrera)

The new server variables are `commit_timestamp_buffers`, `multixact_member_buffers`, `multixact_offset_buffers`, `notify_buffers`, `serializable_buffers`, `subtransaction_buffers`, and `transaction_buffers`. `commit_timestamp_buffers`, `transaction_buffers`, and `subtransaction_buffers` scale up automatically with `shared_buffers`.

E.7.3.1.7. Streaming Replication and Recovery

- Add support for incremental file system backup (Robert Haas, Jakub Wartak, Tomas Vondra)

Incremental backups can be created using `pg_basebackup`'s new `--incremental` option. The new application `pg_combinebackup` allows manipulation of base and incremental file system backups.

- Allow the creation of WAL summarization files (Robert Haas, Nathan Bossart, Hubert Depesz Lubaczewski)

These files record the block numbers that have changed within an `LSN` range and are useful for incremental file system backups. This is controlled by the server variables `summarize_wal` and `wal_summary_keep_time`, and introspected with `pg_available_wal_summaries()`, `pg_wal_summary_contents()`, and `pg_get_wal_summarizer_state()`.

- Add the system identifier to file system `backup manifest` files (Amul Sul)

This helps detect invalid WAL usage.

- Allow connection string value `dbname` to be written when `pg_basebackup` writes connection information to `postgresql.auto.conf` (Vignesh C, Hayato Kuroda)
- Add column `pg_replication_slots.invalidation_reason` to report the reason for invalid slots (Shveta Malik, Bharath Rupireddy)
- Add column `pg_replication_slots.inactive_since` to report slot inactivity duration (Bharath Rupireddy)
- Add function `pg_sync_replication_slots()` to synchronize logical replication slots (Hou Zhijie, Shveta Malik, Ajin Cherian, Peter Eisentraut)
- Add the `failover` property to the `replication protocol` (Hou Zhijie, Shveta Malik)

E.7.3.1.8. Logical Replication

- Add application `pg_createsubscriber` to create a logical replica from a physical standby server (Euler Taveira)
- Have `pg_upgrade` migrate valid logical slots and subscriptions (Hayato Kuroda, Hou Zhijie, Vignesh C, Julien Rouhaud, Shlok Kyal)

This allows logical replication to continue quickly after the upgrade. This only works for old PostgreSQL clusters that are version 17 or later.

- Enable the failover of `logical slots` (Hou Zhijie, Shveta Malik, Ajin Cherian)

This is controlled by an optional fifth argument to `pg_create_logical_replication_slot()`.

- Add server variable `sync_replication_slots` to enable failover logical slot synchronization (Shveta Malik, Hou Zhijie, Peter Smith)

- Add logical replication failover control to `CREATE/ALTER SUBSCRIPTION` (Shveta Malik, Hou Zhijie, Ajin Cherian)
- Allow the application of logical replication changes to use `hash` indexes on the subscriber (Hayato Kuroda)

Previously only `btree` indexes could be used for this purpose.

- Improve `logical decoding` performance in cases where there are many subtransactions (Masahiko Sawada)
- Restart apply workers if subscription owner's superuser privileges are revoked (Vignesh C)

This forces reauthentication.

- Add `flush` option to `pg_logical_emit_message()` (Michael Paquier)

This makes the message durable.

- Allow specification of physical standbys that must be synchronized before they are visible to subscribers (Hou Zhijie, Shveta Malik)

The new server variable is `synchronized_standby_slots`.

- Add worker type column to `pg_stat_subscription` (Peter Smith)

E.7.3.2. Utility Commands

- Add new `COPY` option `ON_ERROR ignore` to discard error rows (Damir Belyalov, Atsushi Torikoshi, Alex Shulgin, Jian He, Yugo Nagata)

The default behavior is `ON_ERROR stop`.

- Add new `COPY` option `LOG_VERBOSITY` which reports `COPY FROM` ignored error rows (Bharath Rupireddy)
- Allow `COPY FROM` to report the number of skipped rows during processing (Atsushi Torikoshi)

This appears in system view column `pg_stat_progress_copy.tuples_skipped`.

- In `COPY FROM`, allow easy specification that all columns should be forced null or not null (Zhang Mingli)
- Allow partitioned tables to have identity columns (Ashutosh Bapat)
- Allow `exclusion constraints` on partitioned tables (Paul A. Jungwirth)

As long as exclusion constraints compare partition key columns for equality, other columns can use exclusion constraint-specific comparisons.

- Add clearer `ALTER TABLE` method to set a column to the default statistics target (Peter Eisentraut)

The new syntax is `ALTER TABLE ... SET STATISTICS DEFAULT`; using `SET STATISTICS -1` is still supported.

- Allow `ALTER TABLE` to change a column's generation expression (Amul Sul)

The syntax is `ALTER TABLE ... ALTER COLUMN ... SET EXPRESSION`.

- Allow specification of `table access methods` on partitioned tables (Justin Pryzby, Soumyadeep Chakraborty, Michael Paquier)
- Add `DEFAULT` setting for `ALTER TABLE .. SET ACCESS METHOD` (Michael Paquier)
- Add support for `event triggers` that fire at connection time (Konstantin Knizhnik, Mikhail Gribkov)
- Add event trigger support for `REINDEX` (Garrett Thornburg, Jian He)

- Allow parenthesized syntax for `CLUSTER` options if a table name is not specified (Nathan Bossart)

E.7.3.2.1. `EXPLAIN`

- Allow `EXPLAIN` to report optimizer memory usage (Ashutosh Bapat)

The option is called `MEMORY`.

- Add `EXPLAIN` option `SERIALIZE` to report the cost of converting data for network transmission (Stepan Rutz, Matthias van de Meent)
- Add local I/O block read/write timing statistics to `EXPLAIN`'s `BUFFERS` output (Nazir Bilal Yavuz)
- Improve `EXPLAIN`'s display of SubPlan nodes and output parameters (Tom Lane, Dean Rasheed)
- Add JIT `deform_counter` details to `EXPLAIN` (Dmitry Dolgov)

E.7.3.3. Data Types

- Allow the `interval` data type to support `+/-infinity` values (Joseph Koshakow, Jian He, Ashutosh Bapat)
- Allow the use of an `ENUM` added via `ALTER TYPE` if the type was created in the same transaction (Tom Lane)

This was previously disallowed.

E.7.3.4. `MERGE`

- Allow `MERGE` to modify updatable views (Dean Rasheed)
- Add `WHEN NOT MATCHED BY SOURCE` to `MERGE` (Dean Rasheed)

`WHEN NOT MATCHED` on target rows was already supported.

- Allow `MERGE` to use the `RETURNING` clause (Dean Rasheed)

The new `RETURNING` function `merge_action()` reports on the DML that generated the row.

E.7.3.5. Functions

- Add function `JSON_TABLE()` to convert `JSON` data to a table representation (Nikita Glukhov, Teodor Sigaev, Oleg Bartunov, Alexander Korotkov, Andrew Dunstan, Amit Langote, Jian He)

This function can be used in the `FROM` clause of `SELECT` queries as a tuple source.

- Add SQL/JSON constructor functions `JSON()`, `JSON_SCALAR()`, and `JSON_SERIALIZE()` (Nikita Glukhov, Teodor Sigaev, Oleg Bartunov, Alexander Korotkov, Andrew Dunstan, Amit Langote)
- Add SQL/JSON query functions `JSON_EXISTS()`, `JSON_QUERY()`, and `JSON_VALUE()` (Nikita Glukhov, Teodor Sigaev, Oleg Bartunov, Alexander Korotkov, Andrew Dunstan, Amit Langote, Peter Eisen-traut, Jian He)
- Add `jsonpath` methods to convert `JSON` values to other `JSON` data types (Jeevan Chalke)

The `jsonpath` methods are `.bigint()`, `.boolean()`, `.date()`, `.decimal([precision [, scale]])`, `.integer()`, `.number()`, `.string()`, `.time()`, `.time_tz()`, `.timestamp()`, and `.timestamp_tz()`.

- Add `to_timestamp()` time zone format specifiers (Tom Lane)

`TZ` accepts time zone abbreviations or numeric offsets, while `OF` accepts only numeric offsets.

- Allow the session `time zone` to be specified by `AT LOCAL` (Vik Fearing)

This is useful when converting adding and removing time zones from time stamps values, rather than specifying the literal session time zone.

- Add functions `uuid_extract_timestamp()` and `uuid_extract_version()` to return UUID information (Andrey Borodin)
- Add functions to generate random numbers in a specified range (Dean Rasheed)

The functions are `random(min, max)` and they take values of type `integer`, `bigint`, and `numeric`.

- Add functions to convert integers to binary and octal strings (Eric Radman, Nathan Bossart)

The functions are `to_bin()` and `to_oct()`.

- Add Unicode informational functions (Jeff Davis)

Function `unicode_version()` returns the Unicode version, `icu_unicode_version()` returns the ICU version, and `unicode_assigned()` returns if the characters are assigned Unicode codepoints.

- Add function `xmltext()` to convert text to a single XML text node (Jim Jones)
- Add function `to_regtypemod()` to return the type modifier of a type specification (David Wheeler, Erik Wienhold)
- Add `pg_basetype()` function to return a domain's base type (Steve Chavez)
- Add function `pg_column_toast_chunk_id()` to return a value's **TOAST** identifier (Yugo Nagata)

This returns `NULL` if the value is not stored in TOAST.

E.7.3.6. PL/pgSQL

- Allow `plpgsql %TYPE` and `%ROWTYPE` specifications to represent arrays of non-array types (Quan Zongliang, Pavel Stehule)
- Allow `plpgsql %TYPE` specification to reference composite column (Tom Lane)

E.7.3.7. libpq

- Add `libpq` function to change role passwords (Joe Conway)

The new function, `PQchangePassword()`, hashes the new password before sending it to the server.

- Add `libpq` functions to close portals and prepared statements (Jelte Fennema-Nio)

The functions are `PQclosePrepared()`, `PQclosePortal()`, `PQsendClosePrepared()`, and `PQsendClosePortal()`.

- Add `libpq` API which allows for blocking and non-blocking **cancel requests**, with encryption if already in use (Jelte Fennema-Nio)

Previously only blocking, unencrypted cancel requests were supported.

- Add `libpq` function `PQsocketPoll()` to allow polling of network sockets (Tristan Partin, Tom Lane)
- Add `libpq` function `PQsendPipelineSync()` to send a pipeline synchronization point (Anton Kirilov)

This is similar to `PQpipelineSync()` but it does not flush to the server unless the size threshold of the output buffer is reached.

- Add `libpq` function `PQsetChunkedRowsMode()` to allow retrieval of results in chunks (Daniel Vérité)
- Allow TLS connections without requiring a network round-trip negotiation (Greg Stark, Heikki Linnakangas, Peter Eisentraut, Michael Paquier, Daniel Gustafsson)

This is enabled with the client-side option `sslnegotiation=direct`, requires ALPN, and only works on PostgreSQL 17 and later servers.

- Improve `psql` display of default and empty privileges (Erik Wienhold, Laurenz Albe)

Command `\dp` now displays `(none)` for empty privileges; default still displays as empty.

- Have backslash commands honor `\pset null` (Erik Wienhold, Laurenz Albe)

Previously `\pset null` was ignored.

- Allow `psql`'s `\watch` to stop after a minimum number of rows returned (Greg Sabino Mullane)

The parameter is `min_rows`.

- Allow `psql` connection attempts to be canceled with control-C (Tristan Partin)
- Allow `psql` to honor `FETCH_COUNT` for non-`SELECT` queries (Daniel Vérité)
- Improve `psql` tab completion (Dagfinn Ilmari Mannsåker, Gilles Darold, Christoph Heiss, Steve Chavez, Vignesh C, Pavel Borisov, Jian He)

E.7.3.9. Server Applications

- Add application `pg_walsummary` to dump WAL summary files (Robert Haas)
- Allow `pg_dump`'s large objects to be restorable in batches (Tom Lane)

This allows the restoration of many large objects to avoid transaction limits and to be restored in parallel.

- Add `pg_dump` option `--exclude-extension` (Ayush Vatsa)
- Allow `pg_dump`, `pg_dumpall`, and `pg_restore` to specify include/exclude objects in a file (Pavel Stehule, Daniel Gustafsson)

The option is called `--filter`.

- Add the `--sync-method` parameter to several client applications (Justin Pryzby, Nathan Bossart)

The applications are `initdb`, `pg_basebackup`, `pg_checksums`, `pg_dump`, `pg_rewind`, and `pg_upgrade`.

- Add `pg_restore` option `--transaction-size` to allow object restores in transaction batches (Tom Lane)

This allows the performance benefits of transaction batches without the problems of excessively large transaction blocks.

- Change `pgbench` debug mode option from `-d` to `--debug` (Greg Sabino Mullane)

Option `-d` is now used for the database name, and the new `--dbname` option can be used as well.

- Add `pgbench` option `--exit-on-abort` to exit after any client aborts (Yugo Nagata)
- Add `pgbench` command `\syncpipeline` to allow sending of sync messages (Anthonin Bonnefoy)
- Allow `pg_archivecleanup` to remove backup history files (Atsushi Torikoshi)

The option is `--clean-backup-history`.

- Add some long options to `pg_archivecleanup` (Atsushi Torikoshi)

The long options are `--debug`, `--dry-run`, and `--strip-extension`.

- Allow `pg_basebackup` and `pg_receivewal` to use `dbname` in their connection specification (Jelte Fennema-Nio)

This is useful for connection poolers that are sensitive to the database name.

- Add `pg_upgrade` option `--copy-file-range` (Thomas Munro)

This is supported on Linux and FreeBSD.

- Allow `reindexdb --index` to process indexes from different tables in parallel (Maxim Orlov, Svetlana Derevyanko, Alexander Korotkov)
- Allow `reindexdb`, `vacuumdb`, and `clusterdb` to process objects in all databases matching a pattern (Nathan Bossart)

The new option `--all` controls this behavior.

E.7.3.10. Source Code

- Remove support for OpenSSL 1.0.1 (Michael Paquier)
- Allow tests to pass in OpenSSL FIPS mode (Peter Eisentraut)
- Use CPU AVX-512 instructions for bit counting (Paul Amonson, Nathan Bossart, Ants Aasma)
- Require LLVM version 10 or later (Thomas Munro)
- Use native CRC instructions on 64-bit LoongArch CPUs (Xudong Yang)
- Remove AIX support (Heikki Linnakangas)
- Remove the Microsoft Visual Studio-specific PostgreSQL build option (Michael Paquier)

Meson is now the only available method for Visual Studio builds.

- Remove configure option `--disable-thread-safety` (Thomas Munro, Heikki Linnakangas)

We now assume all supported platforms have sufficient thread support.

- Remove configure option `--with-CC` (Heikki Linnakangas)

Setting the `CC` environment variable is now the only supported method for specifying the compiler.

- User-defined data type receive functions will no longer receive their data null-terminated (David Rowley)
- Add incremental `JSON` parser for use with huge `JSON` documents (Andrew Dunstan)
- Convert top-level `README` file to Markdown (Nathan Bossart)
- Remove no longer needed top-level `INSTALL` file (Tom Lane)
- Remove make's `distprep` option (Peter Eisentraut)
- Add make support for Android shared libraries (Peter Eisentraut)
- Add backend support for injection points (Michael Paquier)

This is used for server debugging and they must be enabled at server compile time.

- Add dynamic shared memory registry (Nathan Bossart)

This allows shared libraries which are not initialized at startup to coordinate dynamic shared memory access.

- Fix `emit_log_hook` to use the same time value as other log records for the same query (Kambam Vinay, Michael Paquier)
- Improve documentation for using `jsonpath` for predicate checks (David Wheeler)

E.7.3.11. Additional Modules

- Allow joins with non-join qualifications to be pushed down to foreign servers and custom scans (Richard Guo, Etsuro Fujita)

Foreign data wrappers and custom scans will need to be modified to handle these cases.

- Allow pushdown of `EXISTS` and `IN` subqueries to `postgres_fdw` foreign servers (Alexander Pyhalov)

- Increase the default foreign data wrapper tuple cost (David Rowley, Umair Shahid)

This value is used by the optimizer.

- Allow [dblink](#) database operations to be interrupted (Noah Misch)
- Allow the creation of hash indexes on ltree columns (Tommy Pavlicek)

This also enables hash join and hash aggregation on ltree columns.

- Allow unaccent character translation rules to contain whitespace and quotes (Michael Paquier)

The syntax for the `unaccent.rules` file has changed.

- Allow [amcheck](#) to check for unique constraint violations using new option `--checkunique` (Anastasia Lubennikova, Pavel Borisov, Maxim Orlov)
- Allow citext tests to pass in OpenSSL FIPS mode (Peter Eisentraut)
- Allow [pgcrypto](#) tests to pass in OpenSSL FIPS mode (Peter Eisentraut)
- Remove some unused [SPI](#) macros (Bharath Rupireddy)
- Allow [ALTER OPERATOR](#) to set more optimization attributes (Tommy Pavlicek)

This is useful for extensions.

- Allow extensions to define [custom wait events](#) (Masahiro Ikeda)

Custom wait events have been added to [postgres_fdw](#) and [dblink](#).

- Add `pgbuffercache` function `pg_buffercache_evict()` to allow shared buffer eviction (Palak Chaturvedi, Thomas Munro)

This is useful for testing.

E.7.3.11.1. [pg_stat_statements](#)

- Replace [CALL](#) parameters in `pg_stat_statements` with placeholders (Sami Imseih)
- Replace savepoint names stored in `pg_stat_statements` with placeholders (Greg Sabino Mullane)

This greatly reduces the number of entries needed to record [SAVEPOINT](#), [RELEASE SAVEPOINT](#), and [ROLLBACK TO SAVEPOINT](#) commands.

- Replace the two-phase commit GIDs stored in `pg_stat_statements` with placeholders (Michael Paquier)

This greatly reduces the number of entries needed to record [PREPARE TRANSACTION](#), [COMMIT PREPARED](#), and [ROLLBACK PREPARED](#).

- Track [DEALLOCATE](#) in `pg_stat_statements` (Dagfinn Ilmari Mannsåker, Michael Paquier)

`DEALLOCATE` names are stored in `pg_stat_statements` as placeholders.

- Add local I/O block read/write timing statistics columns of `pg_stat_statements` (Nazir Bilal Yavuz)

The new columns are `local_blk_read_time` and `local_blk_write_time`.

- Add JIT `deform_counter` details to `pg_stat_statements` (Dmitry Dolgov)
- Add optional fourth argument (`minmax_only`) to `pg_stat_statements_reset()` to allow for the resetting of only min/max statistics (Andrei Zubkov)

This argument defaults to `false`.

- Add `pg_stat_statements` columns `stats_since` and `minmax_stats_since` to track entry creation time and last min/max reset time (Andrei Zubkov)

E.7.4. Acknowledgments

The following individuals (in alphabetical order) have contributed to this release as patch authors, committers, reviewers, testers, or reporters of issues.

Abhijit Menon-Sen
Adnan Dautovic
Aidar Imamov
Ajin Cherian
Akash Shankaran
Akshat Jaimini
Alaa Attya
Aleksander Alekseev
Aleksej Orlov
Alena Rybakina
Alex Hsieh
Alex Malek
Alex Shulgin
Alex Work
Alexander Korotkov
Alexander Kozhemyakin
Alexander Kuzmenkov
Alexander Lakhin
Alexander Pyhalov
Alexey Palazhchenko
Alfons Kemper
Álvaro Herrera
Amadeo Gallardo
Amit Kapila
Amit Langote
Amul Sul
Anastasia Lubennikova
Anatoly Zaretsky
Andreas Karlsson
Andreas Ulbrich
Andrei Lepikhov
Andrei Zubkov
Andres Freund
Andrew Alsup
Andrew Atkinson
Andrew Bille
Andrew Dunstan
Andrew Kane
Andrey Borodin
Andrey Rachitskiy
Andrey Sokolov
Andy Fan
Anthonin Bonnefoy
Anthony Hsu
Anton Kirilov
Anton Melnikov
Anton Voloshin
Antonin Houska
Ants Aasma
Antti Lampinen
Aramaki Zyake
Artem Anisimov
Artur Zakirov

Ashutosh Bapat
Ashutosh Sharma
Atsushi Torikoshi
Attila Gulyás
Ayush Tiwari
Ayush Vatsa
Bartosz Chrol
Benoît Ryder
Bernd Helmle
Bertrand Drouvot
Bharath Rupireddy
Bo Andreson
Boshomi Phenix
Bowen Shi
Boyu Yang
Bruce Momjian
Cameron Vogt
Cary Huang
Cédric Villemain
Changhong Fei
Chantal Keller
Chapman Flack
Chengxi Sun
Chris Travers
Christian Maurer
Christian Stork
Christoph Berg
Christoph Heiss
Christophe Courtois
Christopher Kline
Claudio Freire
Colin Caine
Corey Huinker
Curt Kolovson
Dag Lem
Dagfinn Ilmari Mannsåker
Damir Belyalov
Daniel Fredouille
Daniel Gustafsson
Daniel Shelepanov
Daniel Vérité
Daniel Westermann
Darren Rush
Dave Cramer
Dave Page
David Christensen
David Cook
David G. Johnston
David Geier
David Hillman
David Perez
David Rowley
David Steele
David Wheeler
David Zhang
Dean Rasheed
Denis Erokhin
Denis Laxalde

Devrim Gündüz
Dilip Kumar
Dimitrios Apostolou
Dmitry Dolgov
Dmitry Koval
Dmitry Vasiliev
Dominique Devienne
Dong Wook Lee
Donghang Lin
Dongming Liu
Drew Callahan
Drew Kimball
Dzmitry Jachnik
Egor Chindyaskin
Egor Rogov
Ekaterina Kiryanova
Elena Indrupskaya
Elizabeth Christensen
Emre Hasegeli
Eric Cyr
Eric Mutta
Eric Radman
Eric Ridge
Erik Rijkers
Erik Wienhold
Erki Eessaar
Ethan Mertz
Etsuro Fujita
Eugen Konkov
Euler Taveira
Evan Macbeth
Evgeny Morozov
Fabien Coelho
Fabrízio de Royes Mello
Farias de Oliveira
Feliphe Pozzer
Fire Emerald
Flavien Guedez
Floris Van Nee
Francesco Degrassi
Frank Streitzig
Gabriele Bartolini
Garrett Thornburg
Gavin Flower
Gavin Panella
Gilles Darold
Gilles Parc
Grant Gryczan
Greg Nancarrow
Greg Sabino Mullane
Greg Stark
Gurjeet Singh
Haiying Tang
Hajime Matsunaga
Hal Takahara
Hanefi Onaldi
Hannu Krosing
Hans Buschmann

Hao Wu
Hao Zhang
Hayato Kuroda
Heikki Linnakangas
Hemanth Sandrana
Himanshu Upadhyaya
Hironobu Suzuki
Holger Reise
Hongxu Ma
Hongyu Song
Horst Reiterer
Hubert Lubaczewski
Hywel Carver
Ian Barwick
Ian Ilyasov
Ilya Nenashev
Isaac Morland
Israel Barth Rubio
Ivan Kartyshov
Ivan Kolombet
Ivan Lazarev
Ivan Panchenko
Ivan Trofimov
Jacob Champion
Jacob Speidel
Jacques Combrink
Jaime Casanova
Jakub Wartak
James Coleman
James Pang
Jani Rahkola
Japin Li
Jeevan Chalke
Jeff Davis
Jeff Janes
Jelte Fennema-Nio
Jeremy Schneider
Jian Guo
Jian He
Jim Jones
Jim Keener
Jim Nasby
Jingtang Zhang
Jingxian Li
Jingzhou Fu
Joe Conway
Joel Jacobson
John Ekins
John Hsu
John Morris
John Naylor
John Russell
Jonathan Katz
Jordi Gutiérrez
Joseph Koshakow
Josh Kupersmidt
Joshua D. Drake
Joshua Uyehara

Jubilee Young
Julien Rouhaud
Junwang Zhao
Justin Pryzby
Kaido Vaikla
Kambam Vinay
Karen Talarico
Karina Litskevich
Karl O. Pinc
Kashif Zeeshan
Kim Johan Andersson
Kirill Reshke
Kirk Parker
Kirk Wolak
Kisoon Kwon
Koen De Groote
Kohei KaiGai
Kong Man
Konstantin Knizhnik
Kouhei Sutou
Krishnakumar R
Kuntal Ghosh
Kurt Roeckx
Kyotaro Horiguchi
Lang Liu
Lars Kanis
Laurenz Albe
Lauri Laanmets
Legs Mansion
Lukas Fittl
Magnus Hagander
Mahendrakar Srinivasarao
Maiquel Grassi
Manos Emmanouilidis
Marcel Hofstetter
Marcos Pegoraro
Marian Krucina
Marina Polyakova
Mark Dilger
Mark Guertin
Mark Sloan
Markus Winand
Marlene Reiterer
Martín Marqués
Martin Nash
Martin Schlossarek
Masahiko Sawada
Masahiro Ikeda
Masaki Kuwamura
Masao Fujii
Mason Sharp
Matheus Alcantara
Mats Kindahl
Matthias Kuhn
Matthias van de Meent
Maxim Boguk
Maxim Orlov
Maxim Yablokov

Maxime Boyer
Melanie Plageman
Melih Mutlu
Merlin Moncure
Micah Gate
Michael Banck
Michael Bondarenko
Michael Paquier
Michael Wang
Michael Zhilin
Michail Nikolaev
Michal Bartak
Michal Kleczek
Mikhail Gribkov
Mingli Zhang
Miroslav Bendik
Mitsuru Hinata
Moaaz Assali
Muralikrishna Bandaru
Nathan Bossart
Nazir Bilal Yavuz
Neil Tiffin
Ngigi Waithaka
Nikhil Benesch
Nikhil Raj
Nikita Glukhov
Nikita Kalinin
Nikita Malakhov
Nikolay Samokhvalov
Nikolay Shaplov
Nisha Moond
Nishant Sharma
Nitin Jadhav
Noah Misch
Noriyoshi Shinoda
Ole Peder Brandtzæg
Oleg Bartunov
Oleg Sibiryakov
Oleg Tselebrovskiy
Olleg Samoylov
Onder Kalaci
Ondrej Navratil
Pablo Khara
Palak Chaturvedi
Pantelis Theodosiou
Paul Amonson
Paul Jungwirth
Pavel Borisov
Pavel Kulakov
Pavel Luzanov
Pavel Stehule
Pavlo Golub
Pedro Gallegos
Pete Storer
Peter Eisentraut
Peter Geoghegan
Peter Smith
Philip Warner

Philipp Salvisberg
Pierre Ducroquet
Pierre Fortin
Przemyslaw Sztoch
Quynh Tran
Raghuveer Devulapalli
Ranier Vilela
Reid Thompson
Rian McGuire
Richard Guo
Richard Vesely
Ridvan Korkmaz
Robert Haas
Robert Scott
Robert Treat
Roberto Mello
Robins Tharakan
Roman Lozko
Ronan Dunklau
Rui Zhao
Ryo Matsumura
Ryoga Yoshida
Sameer Kumar
Sami Imseih
Samuel Dussault
Sanjay Minni
Satoru Koizumi
Sebastian Skalacki
Sergei Glukhov
Sergei Kornilov
Sergey Prokhorenko
Sergey Sargsyan
Sergey Shinderuk
Shaozhong Shi
Shaun Thomas
Shay Rojansky
Shihao Zhong
Shinya Kato
Shlok Kyal
Shruthi Gowda
Shubham Khanna
Shulin Zhou
Shveta Malik
Simon Riggs
Soumyadeep Chakraborty
Sravan Velagandula
Stan Hu
Stepan Neretin
Stepan Rutz
Stéphane Schildknecht
Stephane Tachaires
Stephen Frost
Steve Atkins
Steve Chavez
Suraj Khamkar
Suraj Kharage
Svante Richter
Svetlana Derevyanko

Sylvain Frandaz
Takayuki Tsunakawa
Tatsuo Ishii
Tatsuro Yamada
Tender Wang
Teodor Sigaev
Thom Brown
Thomas Munro
Tim Carey-Smith
Tim Needham
Tim Palmer
Tobias Bussmann
Tom Lane
Tomas Vondra
Tommy Pavlicek
Tomonari Katsumata
Tristan Partin
Tristen Raab
Tung Nguyen
Umair Shahid
Uwe Binder
Valerie Woolard
Vallimaharajan G
Vasya Boytsov
Victor Wagner
Victor Yegorov
Victoria Shepard
Vidushi Gupta
Vignesh C
Vik Fearing
Viktor Leis
Vinayak Pokale
Vitaly Burovoy
Vojtech Benes
Wei Sun
Wei Wang
Wenjiang Zhang
Will Mortensen
Willi Mann
Wolfgang Walther
Xiang Liu
Xiaoran Wang
Xing Guo
Xudong Yang
Yahor Yuzefovich
Yajun Hu
Yaroslav Saburov
Yong Li
Yongtao Huang
Yugo Nagata
Yuhang Qiu
Yuki Seino
Yura Sokolov
Yurii Rashkovskii
Yuuki Fujii
Yuya Watari
Yves Colin
Zhihong Yu

Zhijie Hou
Zongliang Quan
Zubeyr Eryilmaz
Zuming Jiang

E.8. Prior Releases

Release notes for prior versions can be found online. At the time of Postgres Pro Shardman 17 release, these prior versions were supported:

- Postgres Pro Shardman 14: <https://postgrespro.com/docs/shardman/14/release>

Appendix F. Additional Supplied Modules and Extensions Shipped in postgrespro-sdm-17-contrib

This appendix, [Appendix G](#), and [Appendix H](#) contain information on the optional components available in the Postgres Pro Shardman distribution. These include porting tools, analysis utilities, and plug-in features that are not part of the core Postgres Pro system. They are separate mainly because they address a limited audience or are too experimental to be part of the main source tree. This does not preclude their usefulness.

This appendix and [Appendix G](#) cover the extensions and other server plug-in modules. [Appendix H](#) covers the utility programs.

Many components supply new user-defined functions, operators, or types, packaged as *extensions*. To make use of one of these extensions, after you have installed the code you need to register the new SQL objects in the database system. This is done by executing a [CREATE EXTENSION](#) command. In a fresh database, you can simply do

```
CREATE EXTENSION extension_name;
```

This command registers the new SQL objects in the current database only, so you need to run it in every database in which you want the extension's facilities to be available. Alternatively, run it in database `template1` so that the extension will be copied into subsequently-created databases by default.

For all extensions, the `CREATE EXTENSION` command must be run by a database superuser, unless the extension is considered “trusted”. Trusted extensions can be run by any user who has `CREATE` privilege on the current database. Extensions that are trusted are identified as such in the sections that follow. Generally, trusted extensions are ones that cannot provide access to outside-the-database functionality.

The following extensions are trusted in a default installation:

[pgcrypto](#) [uuid-oss](#)

Many extensions allow you to install their objects in a schema of your choice. To do that, add `SCHEMA schema_name` to the `CREATE EXTENSION` command. By default, the objects will be placed in your current creation target schema, which in turn defaults to `public`.

Note, however, that some of these components are not “extensions” in this sense, but are loaded into the server in some other way, for instance by way of [shared_preload_libraries](#). See the documentation of each component for details.

This appendix contains modules and extensions that are made available in Postgres Pro Shardman as a separate subpackage `postgrespro-sdm-17-contrib`. Note that the `toastapi` extension to add custom TOASTers is also available in `postgrespro-sdm-17-contrib` (see [Section 66.2.3](#) for details).

F.1. amcheck — tools to verify table and index consistency

The `amcheck` module provides functions that allow you to verify the logical consistency of the structure of relations.

The B-Tree checking functions verify various *invariants* in the structure of the representation of particular relations. The correctness of the access method functions behind index scans and other important operations relies on these invariants always holding. For example, certain functions verify, among other things, that all B-Tree pages have items in “logical” order (e.g., for B-Tree indexes on `text`, index tuples should be in collated lexical order). If that particular invariant somehow fails to hold, we can expect binary searches on the affected page to incorrectly guide index scans, resulting in wrong answers to SQL queries. If the structure appears to be valid, no error is raised. While these checking functions are run, the `search_path` is temporarily changed to `pg_catalog`, `pg_temp`.

Verification is performed using the same procedures as those used by index scans themselves, which may be user-defined operator class code. For example, B-Tree index verification relies on comparisons made with one or more B-Tree support function 1 routines. See [Section 39.16.3](#) for details of operator class support functions.

Unlike the B-Tree checking functions which report corruption by raising errors, the heap checking function `verify_heapam` checks a table and attempts to return a set of rows, one row per corruption detected. Despite this, if facilities that `verify_heapam` relies upon are themselves corrupted, the function may be unable to continue and may instead raise an error.

Permission to execute `amcheck` functions may be granted to non-superusers, but before granting such permissions careful consideration should be given to data security and privacy concerns. Although the corruption reports generated by these functions do not focus on the contents of the corrupted data so much as on the structure of that data and the nature of the corruptions found, an attacker who gains permission to execute these functions, particularly if the attacker can also induce corruption, might be able to infer something of the data itself from such messages.

Note

Using `amcheck 1.1.1` is not recommended as the `bt_index_parent_check` function works incorrectly in that version.

F.1.1. Functions

`bt_index_check(index regclass, heapallindexed boolean, checkunique boolean)` returns void

`bt_index_check` tests that its target, a B-Tree index, respects a variety of invariants. Example usage:

```
test=# SELECT bt_index_check(index => c.oid, heapallindexed => i.indisunique),
        c.relname,
        c.relpages
FROM pg_index i
JOIN pg_opclass op ON i.indclass[0] = op.oid
JOIN pg_am am ON op.opcmethod = am.oid
JOIN pg_class c ON i.indexrelid = c.oid
JOIN pg_namespace n ON c.relnamespace = n.oid
WHERE am.amname = 'btree' AND n.nspname = 'pg_catalog'
-- Don't check temp tables, which may be from another session:
AND c.relpersistence != 't'
-- Function may throw an error when this is omitted:
AND c.relkind = 'i' AND i.indisready AND i.indisvalid
ORDER BY c.relpages DESC LIMIT 10;
bt_index_check |                relname                | relpages
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

	pg_depend_reference_index		43
	pg_depend_depender_index		40
	pg_proc_proname_args_nsp_index		31
	pg_description_o_c_o_index		21
	pg_attribute_relid_attnam_index		14
	pg_proc_oid_index		10
	pg_attribute_relid_attnum_index		9
	pg_amproc_fam_proc_index		5
	pg_amop_opr_fam_index		5
	pg_amop_fam_strat_index		5

(10 rows)

This example shows a session that performs verification of the 10 largest catalog indexes in the database “test”. Verification of the presence of heap tuples as index tuples is requested for the subset that are unique indexes. Since no error is raised, all indexes tested appear to be logically consistent. Naturally, this query could easily be changed to call `bt_index_check` for every index in the database where verification is supported.

`bt_index_check` acquires an `AccessShareLock` on the target index and the heap relation it belongs to. This lock mode is the same lock mode acquired on relations by simple `SELECT` statements. `bt_index_check` does not verify invariants that span child/parent relationships, but will verify the presence of all heap tuples as index tuples within the index when *heapallindexed* is `true`. When *checkunique* is `true` `bt_index_check` will check that no more than one among duplicate entries in unique index is visible. When a routine, lightweight test for corruption is required in a live production environment, using `bt_index_check` often provides the best trade-off between thoroughness of verification and limiting the impact on application performance and availability.

`bt_index_parent_check(index regclass, heapallindexed boolean, rootdescend boolean, checkunique boolean)` returns void

`bt_index_parent_check` tests that its target, a B-Tree index, respects a variety of invariants. Optionally, when the *heapallindexed* argument is `true`, the function verifies the presence of all heap tuples that should be found within the index. When *checkunique* is `true` `bt_index_parent_check` will check that no more than one among duplicate entries in unique index is visible. When the optional *rootdescend* argument is `true`, verification re-finds tuples on the leaf level by performing a new search from the root page for each tuple. The checks that can be performed by `bt_index_parent_check` are a superset of the checks that can be performed by `bt_index_check`. `bt_index_parent_check` can be thought of as a more thorough variant of `bt_index_check`: unlike `bt_index_check`, `bt_index_parent_check` also checks invariants that span parent/child relationships, including checking that there are no missing downlinks in the index structure. `bt_index_parent_check` follows the general convention of raising an error if it finds a logical inconsistency or other problem.

A `ShareLock` is required on the target index by `bt_index_parent_check` (a `ShareLock` is also acquired on the heap relation). These locks prevent concurrent data modification from `INSERT`, `UPDATE`, and `DELETE` commands. The locks also prevent the underlying relation from being concurrently processed by `VACUUM`, as well as all other utility commands. Note that the function holds locks only while running, not for the entire transaction.

`bt_index_parent_check`'s additional verification is more likely to detect various pathological cases. These cases may involve an incorrectly implemented B-Tree operator class used by the index that is checked, or, hypothetically, undiscovered bugs in the underlying B-Tree index access method code. Note that `bt_index_parent_check` cannot be used when hot standby mode is enabled (i.e., on read-only physical replicas), unlike `bt_index_check`.

Tip

`bt_index_check` and `bt_index_parent_check` both output log messages about the verification process at `DEBUG1` and `DEBUG2` severity levels. These messages provide detailed information about

the verification process that may be of interest to PostgreSQL developers. Advanced users may also find this information helpful, since it provides additional context should verification actually detect an inconsistency. Running:

```
SET client_min_messages = DEBUG1;
```

in an interactive psql session before running a verification query will display messages about the progress of verification with a manageable level of detail.

```
verify_heapam(relation regclass, on_error_stop boolean, check_toast boolean, skip text,  
startblock bigint, endblock bigint, blkno OUT bigint, offnum OUT integer, attnum OUT  
integer, msg OUT text) returns setof record
```

Checks a table, sequence, or materialized view for structural corruption, where pages in the relation contain data that is invalidly formatted, and for logical corruption, where pages are structurally valid but inconsistent with the rest of the database cluster.

The following optional arguments are recognized:

`on_error_stop`

If true, corruption checking stops at the end of the first block in which any corruptions are found.

Defaults to false.

`check_toast`

If true, toasted values are checked against the target relation's TOAST table.

This option is known to be slow. Also, if the toast table or its index is corrupt, checking it against toast values could conceivably crash the server, although in many cases this would just produce an error.

Defaults to false.

`skip`

If not `none`, corruption checking skips blocks that are marked as all-visible or all-frozen, as specified. Valid options are `all-visible`, `all-frozen` and `none`.

Defaults to `none`.

`startblock`

If specified, corruption checking begins at the specified block, skipping all previous blocks. It is an error to specify a `startblock` outside the range of blocks in the target table.

By default, checking begins at the first block.

`endblock`

If specified, corruption checking ends at the specified block, skipping all remaining blocks. It is an error to specify an `endblock` outside the range of blocks in the target table.

By default, all blocks are checked.

For each corruption detected, `verify_heapam` returns a row with the following columns:

`blkno`

The number of the block containing the corrupt page.

`offnum`

The OffsetNumber of the corrupt tuple.

`attnum`

The attribute number of the corrupt column in the tuple, if the corruption is specific to a column and not the tuple as a whole.

`msg`

A message describing the problem detected.

F.1.2. Optional *heapallindexed* Verification

When the *heapallindexed* argument to B-Tree verification functions is `true`, an additional phase of verification is performed against the table associated with the target index relation. This consists of a “dummy” `CREATE INDEX` operation, which checks for the presence of all hypothetical new index tuples against a temporary, in-memory summarizing structure (this is built when needed during the basic first phase of verification). The summarizing structure “fingerprints” every tuple found within the target index. The high level principle behind *heapallindexed* verification is that a new index that is equivalent to the existing, target index must only have entries that can be found in the existing structure.

The additional *heapallindexed* phase adds significant overhead: verification will typically take several times longer. However, there is no change to the relation-level locks acquired when *heapallindexed* verification is performed.

The summarizing structure is bound in size by `maintenance_work_mem`. In order to ensure that there is no more than a 2% probability of failure to detect an inconsistency for each heap tuple that should be represented in the index, approximately 2 bytes of memory are needed per tuple. As less memory is made available per tuple, the probability of missing an inconsistency slowly increases. This approach limits the overhead of verification significantly, while only slightly reducing the probability of detecting a problem, especially for installations where verification is treated as a routine maintenance task. Any single absent or malformed tuple has a new opportunity to be detected with each new verification attempt.

F.1.3. Using *amcheck* Effectively

amcheck can be effective at detecting various types of failure modes that [data checksums](#) will fail to catch. These include:

- Structural inconsistencies caused by incorrect operator class implementations.

This includes issues caused by the comparison rules of operating system collations changing. Comparisons of datums of a collatable type like `text` must be immutable (just as all comparisons used for B-Tree index scans must be immutable), which implies that operating system collation rules must never change. Though rare, updates to operating system collation rules can cause these issues. More commonly, an inconsistency in the collation order between a primary server and a standby server is implicated, possibly because the *major* operating system version in use is inconsistent. Such inconsistencies will generally only arise on standby servers, and so can generally only be detected on standby servers.

If a problem like this arises, it may not affect each individual index that is ordered using an affected collation, simply because *indexed* values might happen to have the same absolute ordering regardless of the behavioral inconsistency. See [Section 23.1](#) and [Section 23.2](#) for further details about how Postgres Pro uses operating system locales and collations.

- Structural inconsistencies between indexes and the heap relations that are indexed (when *heapallindexed* verification is performed).

There is no cross-checking of indexes against their heap relation during normal operation. Symptoms of heap corruption can be subtle.

- Corruption caused by hypothetical undiscovered bugs in the underlying Postgres Pro access method code, sort code, or transaction management code.

Automatic verification of the structural integrity of indexes plays a role in the general testing of new or proposed Postgres Pro features that could plausibly allow a logical inconsistency to be in-

troduced. Verification of table structure and associated visibility and transaction status information plays a similar role. One obvious testing strategy is to call `amcheck` functions continuously when running the standard regression tests.

- File system or storage subsystem faults where checksums happen to simply not be enabled.

Note that `amcheck` examines a page as represented in some shared memory buffer at the time of verification if there is only a shared buffer hit when accessing the block. Consequently, `amcheck` does not necessarily examine data read from the file system at the time of verification. Note that when checksums are enabled, `amcheck` may raise an error due to a checksum failure when a corrupt block is read into a buffer.

- Corruption caused by faulty RAM, or the broader memory subsystem.

Postgres Pro does not protect against correctable memory errors and it is assumed you will operate using RAM that uses industry standard Error Correcting Codes (ECC) or better protection. However, ECC memory is typically only immune to single-bit errors, and should not be assumed to provide *absolute* protection against failures that result in memory corruption.

When `heapallindexed` verification is performed, there is generally a greatly increased chance of detecting single-bit errors, since strict binary equality is tested, and the indexed attributes within the heap are tested.

Structural corruption can happen due to faulty storage hardware, or relation files being overwritten or modified by unrelated software. This kind of corruption can also be detected with [data page checksums](#).

Relation pages which are correctly formatted, internally consistent, and correct relative to their own internal checksums may still contain logical corruption. As such, this kind of corruption cannot be detected with checksums. Examples include toasted values in the main table which lack a corresponding entry in the toast table, and tuples in the main table with a Transaction ID that is older than the oldest valid Transaction ID in the database or cluster.

Multiple causes of logical corruption have been observed in production systems, including bugs in the PostgreSQL server software, faulty and ill-conceived backup and restore tools, and user error.

Corrupt relations are most concerning in live production environments, precisely the same environments where high risk activities are least welcome. For this reason, `verify_heapam` has been designed to diagnose corruption without undue risk. It cannot guard against all causes of backend crashes, as even executing the calling query could be unsafe on a badly corrupted system. Access to [catalog tables](#) is performed and could be problematic if the catalogs themselves are corrupted.

In general, `amcheck` can only prove the presence of corruption; it cannot prove its absence.

F.1.4. Repairing Corruption

No error concerning corruption raised by `amcheck` should ever be a false positive. `amcheck` raises errors in the event of conditions that, by definition, should never happen, and so careful analysis of `amcheck` errors is often required.

There is no general method of repairing problems that `amcheck` detects. An explanation for the root cause of an invariant violation should be sought. A `REINDEX` may not be effective in repairing corruption.

F.2. aqo — cost-based query optimization

The aqo module is a Postgres Pro Shardman extension for cost-based query optimization. Using machine learning methods, more precisely, a modification of the k-NN algorithm, aqo improves cardinality estimation, which can optimize execution plans and, consequently, speed up query execution.

F.2.1. Description

The aqo module can collect statistics on all the executed queries, excluding the queries that access system relations. The collected statistics is classified by query class. If the queries differ in their constants only, they belong to the same class. For each query class, aqo stores the cardinality quality, planning time, execution time, and execution statistics for machine learning. Based on this data, aqo builds a new query plan and uses it for the next query of the same class. aqo test runs have shown significant performance improvements for complex queries.

aqo saves all the learning data ([aqo_data](#)), queries ([aqo_query_texts](#)), query settings ([aqo_queries](#)), and query execution statistics ([aqo_query_stat](#)) to files. When aqo starts, it loads this data to shared memory. You can access aqo data through functions and views.

When [aqo.advanced](#) is on, and aqo is run in the `intelligent` or `learn` mode, a unique hash value, which is computed from the query tree, is assigned to each query class to identify it and separate the collected statistics. If [aqo.advanced](#) is off, the statistics for all untracked query classes is stored in a common query class with hash 0.

Each query class has an associated separate space, called *feature space*, in which the statistics for this query class is collected. This feature space is identified by a hash value (`fs`), called a *basic hash*, which is the same for queries that only differ in table names, so the learning data is aggregated for such queries. Each feature space has associated *feature subspaces*, where the information about selectivity and cardinality for each query plan node is collected. Each subspace is also identified by a hash value (`fss`).

Query-specific optimization settings are stored in the [aqo_queries](#) view.

F.2.1.1. Limitations

aqo currently has the following limitations:

- Query optimization using the aqo module is not supported on standby.
- Query optimization with aqo does not work with queries that only have temporary objects.
- Query optimization with aqo does not work for queries that contain `IMMUTABLE` functions.
- aqo does not collect statistics on replicas because replicas are read-only. However, aqo may use query execution statistics from the primary if the replica is physical.
- `learn` and `intelligent` modes are not supposed to work for a whole cluster with queries having a dynamically generated structure because these modes store all query class IDs, which are different for all queries in such a workload. Dynamically generated constants are supported, however.

F.2.2. Installation and Setup

The aqo extension is included into Postgres Pro Shardman. Once you have Postgres Pro Shardman installed, complete the following steps to enable aqo:

1. Add aqo to the [shared_preload_libraries](#) parameter in the `postgresql.conf` file:

```
shared_preload_libraries = 'aqo'
```

The aqo library must be preloaded at the server startup, since adaptive query optimization needs to be enabled per cluster.

2. Create the aqo extension using the following query:

```
CREATE EXTENSION aqo;
```

Once the extension is created, you can start optimizing queries.

To disable aqo in the current database, run:

```
DROP EXTENSION aqo;
```

To disable aqo at the cluster level, run:

```
ALTER SYSTEM SET aqo.enable = off;  
SELECT pg_reload_conf();
```

To remove all the aqo data including the collected statistics, call `aqo_reset()`: to remove the data from the current database, run:

```
SELECT aqo_reset();
```

to remove all the data from the aqo storage, run:

```
SELECT aqo_reset(NULL);
```

If you do not want aqo to be loaded at the server restart, remove the line

```
shared_preload_libraries = 'aqo'
```

from the `postgresql.conf` file.

Important

For smooth physical replication transferring aqo data from the primary to a replica, ensure that the same aqo versions are installed on both. You can have different aqo versions installed, but in this case, set [aqo.wal_rw](#) to `off` on both and anticipate no replication.

F.2.2.1. Configuration

aqo behavior mainly depends on the [aqo.enable](#), [aqo.mode](#) and [aqo.advanced](#) configuration parameters. Their default values allow you to start learning in the [aqo basic mode](#) once you just set `aqo.enable` to `on`.

To dynamically change any of these parameters, for example, `mode`, in your current session, run the following command:

```
SET aqo.mode = 'mode';
```

where *mode* is the name of the operation mode to use.

F.2.3. Usage

F.2.3.1. Using aqo in a Basic Mode

By default, [aqo.advanced](#) is off. This sets a recommended, basic, mode, where statistics is collected for plan nodes (identified by `fss`), and the collected machine learning data is used to correct the cardinality error for all queries whose plan contains a certain plan node. Once you set [aqo.enable](#) to `on`, aqo starts learning. Execute queries that you need to optimize several times until the plan is good enough and change [aqo.mode](#) to `frozen`. To apply the machine learning data at the level of the server instance, run the following command:

```
ALTER SYSTEM SET aqo.mode = frozen;  
ALTER SYSTEM SET aqo.enable = on;  
SELECT pg_reload_conf();
```

The machine learning data will be applied not only to the queries on which aqo learned, but to all the queries whose plan contains the nodes for which the statistics was collected. For the machine learning data not affect other queries, set [aqo.advanced](#) to `on` and collect statistics for individual queries. See [Section F.2.3.2](#) for details.

F.2.3.2. Choosing the Operation Mode for Advanced Query Optimization

If you often run queries of the same class, for example, your application limits the number of possible query classes, you can enable `aqo.advanced` and use the `intelligent` mode to improve planning for these queries. In this mode, aqo analyzes each query execution and stores statistics. Statistics on queries of different classes is stored separately. If performance is not improved after 50 iterations, the aqo extension falls back to the default query planner.

Note

You can view the current query plan using the standard Postgres Pro `EXPLAIN` command with the `ANALYZE` option. For details, see the [Section 14.1](#).

Since the `intelligent` mode tries to learn separately for different query classes, aqo may fail to provide performance improvements if queries in the workload are of multiple different classes or if the classes of the queries in the workload are constantly changing. For such workloads, reset the aqo extension to the `controlled` mode, or try to turn off `aqo.advanced`.

When `aqo.advanced` is on, in the `controlled` mode, aqo does not collect statistics for new query classes, so they will not be optimized, but for known query classes, aqo will continue collecting statistics and using optimized planning algorithms. So use the `controlled` mode only after aqo learned in the `learn` or `intelligent` mode.

After aqo has already learned, `controlled` is the mode recommended for production use. To make aqo run in this mode on your whole production cluster, run

```
ALTER SYSTEM SET aqo.mode = 'controlled';  
SELECT pg_reload_conf();
```

When `aqo.advanced` is on, the `learn` mode collects statistics from all the executed queries and updates the data for query classes. This mode is similar to the `intelligent` mode, except that it does not provide intelligent tuning. This mode is not recommended to be used permanently for a whole cluster because it tries aqo optimizations for every query class, even for those that do not need it, and this may lead to an unnecessary computational overhead and cause performance degradation.

Use the `learn` mode with `aqo.advanced` turned off to handle workloads with dynamically generated query structures. Overall performance improvement is not guaranteed. As this mode lacks intelligent tuning, the performance for some queries may even decrease, but this mode is good for a dynamic workload and consumes less memory than the `intelligent` mode.

When `aqo.advanced` is on, in the `auto` mode, aqo behaves like in the `learn` mode, but the `auto` mode prevents the memory overflow by implementing the least recently used (LRU) cache mechanism. This mechanism ensures that the `aqo.dsm_size_max` and the `aqo.fss_max_items` limits are not exceeded when incoming queries are saved. When the `aqo_query_texts_update` function is called in this mode, aqo does not save query texts and does not free the dynamic shared memory used.

For data that may change significantly between queries, enabling the `aqo.delta_rows` parameter may be useful. With this setting, aqo makes predictions based on the planner's row count estimates. For example, if data is deleted from a table, the planner estimates fewer rows than before. aqo can then use this updated information to predict a lower row count without requiring additional learning steps.

If you want to reduce the impact of aqo on query planning and execution, you can use it in the `frozen` mode. In this mode, aqo only reads the collected statistics, but does not collect any new data.

F.2.3.3. Fine-Tuning aqo

You must have superuser rights to access aqo views and configure advanced query settings.

You can view all the processed query classes and their corresponding hash values in the [aqo_query_texts](#) view:

```
SELECT * FROM aqo_query_texts;
```

To find out the class, that is, hash, of a query and aqo mode, enable [aqo.show_hash](#) (boolean) and [aqo.show_details](#) (boolean) environment variables and execute the query. The output will contain something like this:

```
...
Planning Time: 23.538 ms
...
Execution Time: 249813.875 ms
...
Using aqo: true
AQO mode: LEARN
AQO advanced: OFF
...
Query hash: -2439501042637610315
```

Each query class has its own optimization settings. These settings are shown in the [aqo_queries](#) view:

```
SELECT * FROM aqo_queries;
```

You can manually change these settings to adjust the optimization for a particular query class. For example:

```
-- Add a new query class to the aqo_queries view:

SET aqo.enable='on';
SET aqo.advanced='on';
SET aqo.mode='intelligent';
SELECT * FROM a, b WHERE a.id=b.id;
SET aqo.mode='controlled';

-- Disable auto_tuning, enable both learn_aqo and use_aqo
-- for this query class:

SELECT count(*) FROM aqo_queries,
  LATERAL aqo_queries_update(fs, NULL, true, true, false)
WHERE fs = (SELECT fs FROM aqo_query_texts
WHERE query_text LIKE 'SELECT * FROM a, b WHERE a.id=b.id;');

-- Run EXPLAIN ANALYZE while the plan changes:

EXPLAIN ANALYZE SELECT * FROM a, b WHERE a.id=b.id;
EXPLAIN ANALYZE SELECT * FROM a, b WHERE a.id=b.id;

-- Disable learning to stop statistics collection
-- and use the optimized plan:
```

```
SELECT count(*) FROM aqo_queries,
  LATERAL aqo_queries_update(fs, NULL, false, true, false)
WHERE fs = (SELECT fs FROM aqo_query_texts
WHERE query_text LIKE 'SELECT * FROM a, b WHERE a.id=b.id;');
```

To stop intelligent tuning for a particular query class, disable the `auto_tuning` setting:

```
SELECT count(*) FROM aqo_queries,
  LATERAL aqo_queries_update(fs, NULL, NULL, NULL, false)
```

```
WHERE fs = 'hash';
```

where *hash* is the hash value for this query class. As a result, aqo disables automatic change of the `learn_aqo` and `use_aqo` settings.

To disable further learning for a particular query class, use the following command:

```
SELECT count(*) FROM aqo_queries,  
       LATERAL aqo_queries_update(fs, NULL, false, NULL, false)  
       WHERE fs = 'hash';
```

where *hash* is the hash value for this query class.

To fully disable aqo for all queries and use the default Postgres Pro query planner, run:

```
SELECT count(*) FROM aqo_queries,  
       LATERAL aqo_disable_class(fs, NULL)  
       WHERE fs <> 0;
```

To disable aqo for all queries temporarily in the current session or for the whole cluster, but not remove or change the collected statistics and settings, turn off [aqo.enable](#):

```
SET aqo.enable = 'off';
```

or

```
ALTER SYSTEM SET aqo.enable = 'off'
```

F.2.3.4. Sandbox Mode

You can experiment with aqo without touching its main knowledge base. To do this, execute the command

```
SET aqo.sandbox = ON;
```

This turns on the *sandbox mode*, which means that aqo will work in the isolated environment. However, if you turn on [aqo.sandbox](#) in different SQL sessions, they will use the same data.

Data obtained in the sandbox mode does not get replicated. But the sandbox mode can be used on a standby. Moreover, the only way to train aqo on a standby is turning on the sandbox mode when the replication is turned on, that is, `aqo.wal_rw` is true. Without the sandbox mode, aqo will work on the standby as if `aqo.mode = FROZEN`, that is, it will be able to use the existing knowledge base, but not update or extend it.

F.2.4. Reference

F.2.4.1. Configuration Parameters

`aqo.enable` (boolean)

Defines the state of aqo. If set to `off` aqo does not work except when [aqo.force_collect_stat](#) = `on`.

Default: `off`.

`aqo.mode` (text)

Sets the aqo operation mode. Defines how aqo handles new queries. Possible values:

- `intelligent` — saves new queries with `auto_tuning` enabled. See the description of the [aqo_queries](#) view for more details. May disable aqo for a query in the case of average performance reduction. Only works in this way if [aqo.advanced](#) = `on`, otherwise, this mode works exactly like `learn`.
- `learn` — collects statistics on all the executed queries, learns and makes predictions based on these statistics.
- `controlled` — only learns and makes predictions for known queries.

- `frozen` — makes predictions for known queries, but does not learn from any queries.
- `auto` — behaves like `learn` but prevents memory overflow by implementing the LRU cache mechanism.

Default: `learn`.

`aqo.advanced` (boolean)

Enables the advanced learning routine, which saves separate learning statistics for each query class. Also allows fine-tuning the `use_aqo` and `learn_aqo` settings in the [aqo_queries](#) view. Fine-tuned query settings in the `aqo_query` view continue to work if `aqo.advanced` is disabled.

Default: `off`.

`aqo.force_collect_stat` (boolean)

Collects statistics on query executions in all aqo modes and even if `aqo.enable` is `off`.

Default: `off`.

`aqo.show_details` (boolean)

Adds some details to `EXPLAIN` output of a query, such as the prediction or feature-subspace hash, and shows some additional aqo-specific on-screen information.

Default: `on`.

`aqo.show_hash` (boolean)

Shows a hash value that uniquely identifies the class of queries or class of plan nodes. aqo uses the native query ID to identify a query class for consistency with other extensions, such as [pg_stat_statements](#). So, the query ID can be taken from the `Query hash` field in `EXPLAIN ANALYZE` output of a query.

Default: `on`.

`aqo.join_threshold` (integer)

Ignores queries that contain smaller number of joins, which means that statistics for such queries is not collected.

Default: 0 (no queries are ignored).

`aqo.learn_statement_timeout` (boolean)

Learns on a plan interrupted by the statement timeout.

Default: `off`.

`aqo.statement_timeout` (integer)

Defines the initial value of the *smart statement timeout*, in milliseconds, which is needed to limit the execution time when manually training aqo on special queries with a poor cardinality forecast. aqo can dynamically change the value of the smart statement timeout during this training. When the cardinality estimation error on nodes exceeds 0.1, the value of `aqo.statement_timeout` is automatically incremented exponentially, but remains not greater than [statement_timeout](#).

Default: 0.

`aqo.wide_search` (boolean)

Enables searching neighbors with the same feature subspace among different query classes. Only has an effect if [aqo.advanced](#) = `on`.

Default: `off`.

`aqo.min_neighbors_for_predicting` (integer)

Defines how many samples collected in previous executions of the query will be used to predict the cardinality next time. If there are fewer of them, aqo will not make any prediction. A too large value may affect performance, but a too small value may reduce the prediction quality.

Default: 3.

`aqo.predict_with_few_neighbors` (boolean)

Enables aqo to make predictions with fewer neighbors than specified by [aqo.min_neighbors_for_predicting](#). When set to `off`, then aqo learns, but does not make predictions until the execution count for the query with different constants reaches 3 (default for `aqo.min_neighbors_for_predicting`).

Default: `on`.

`aqo.fs_max_items` (integer)

Defines the maximum number of feature spaces that aqo can operate with. When this number is exceeded, learning on new query classes will no longer occur, and they will not appear in the views accordingly. This parameter can only be set at server start.

Default: 10000.

`aqo.fss_max_items` (integer)

Defines the maximum number of feature subspaces that aqo can operate with. When this number is exceeded, the selectivity and cardinality for new query plan nodes will no longer be collected, and new feature subspaces will not appear in the [aqo_data](#) view accordingly. This parameter can only be set at server start.

Default: 100000.

`aqo.querytext_max_size` (integer)

Defines the maximum size of the query in the [aqo_query_texts](#) view. This parameter can only be set at server start.

Default: 1000.

`aqo.dsm_size_max` (integer)

Defines the maximum size of dynamic shared memory, in MB, that aqo can allocate to store learning data and query texts. If set to a number that is less than the size of the saved aqo data, the server will not start. This parameter can only be set at server start.

Default: 100.

`aqo.wal_rw` (boolean)

Enables physical replication and allows complete aqo data recovery after failure. When set to `off` on the primary, no data is transferred from it to a replica. When set to `off` on a replica, any data transferred from the primary is ignored. With this value, when the server fails, data can only be restored as of the last checkpoint. This parameter can only be set at server start.

Default: `on`.

`aqo.sandbox` (boolean)

Enables reserving a separate memory area in shared memory to be used by a primary or standby node, which allows collecting and using statistics with the data in this memory area. If enabled on the

primary, the extension uses the separate shared memory area that is not replicated to the standby. Changing the value of this parameter resets the aqo cache. Only superusers can change this setting.

Default: `off`.

`aqo.delta_rows` (boolean)

Enables a learning mechanism where aqo adjusts the planner's row count estimates with its own predictions. If disabled, aqo uses its own predictions.

Default: `off`.

F.2.4.2. Views

F.2.4.2.1. `aqo_query_texts`

The `aqo_query_texts` view classifies all the query classes processed by aqo. For each query class, the view shows the text of the first analyzed query of this class.

Table F.1. `aqo_query_texts` View

Column Name	Description
<code>queryid</code>	The unique identifier of the query.
<code>dbid</code>	The identifier of the database.
<code>fs</code>	The identifier of the feature space.
<code>query_text</code>	Text of the first analyzed query of the given class. The query text length is limited by <code>aqo.querytext_max_size</code> .

F.2.4.2.2. `aqo_queries`

The `aqo_queries` view shows optimization settings for different query classes. One query executed in two different databases is stored twice although the `fs` is the same.

Table F.2. `aqo_queries` View

Setting	Description
<code>fs</code>	The identifier of the feature space.
<code>dbid</code>	The identifier of the database in which the query was executed.
<code>learn_aqo</code>	Shows whether statistics collection for this query class is enabled.
<code>use_aqo</code>	Shows whether the aqo cardinality prediction for the next execution of this query class is enabled.
<code>auto_tuning</code>	<p>Shows whether aqo can dynamically change <code>use_aqo</code> and <code>learn_aqo</code> settings for this query class. By default, set to <code>true</code> for new queries if <code>aqo.advanced</code> is on and <code>aqo.mode = intelligent</code>.</p> <p>When <code>auto_tuning</code> is on, if for several successive executions of a query for which <code>use_aqo</code> is off, the cardinality error remains sufficiently small and stable, aqo turns on <code>use_aqo</code>.</p> <p>For queries with <code>learn_aqo = true</code> (it is so for new queries), several first executions are done both using aqo and without it. The faster the query is executed compared to the execution with</p>

Setting	Description
	the standard planner, the more likely aqo will be used for the next query execution. If after a certain number of executions the execution time with aqo appears to be worse than with the standard planner, aqo will never be used for this query class: <code>auto_tuning</code> , <code>use_aqo</code> and <code>learn_aqo</code> are set to <code>off</code> .
<code>smart_timeout</code>	The value of the smart statement timeout for this query class. The initial value of the smart statement timeout for any query is defined by the statement_timeout configuration parameter.
<code>count_increase_timeout</code>	Shows how many times the smart statement timeout increased for this query class.

F.2.4.2.3. aqo_data

The `aqo_data` view shows machine learning data for cardinality estimation refinement. The number of rows is limited by [aqo.fss_max_items](#). To discard all the collected statistics for a particular query class, you can delete all rows from `aqo_data` with the corresponding `fs`.

Table F.3. aqo_data View

Data	Description
<code>fs</code>	The identifier (hash) of the feature space.
<code>fss</code>	The identifier (hash) of the feature subspace.
<code>dbid</code>	The identifier of the database.
<code>delta_rows</code>	If true, aqo makes predictions based on the planner's estimates, otherwise <code>false</code> .
<code>nfeatures</code>	Feature-subspace size for the query plan node.
<code>features</code>	Logarithm of the selectivity which the cardinality prediction is based on.
<code>targets</code>	Cardinality logarithm for the query plan node.
<code>reliability</code>	Confidence level of the learning statistics. Equals: <ul style="list-style-type: none"> • 1 (default) — indicates data obtained after normal execution of a query • 0.1 — indicates data obtained from a partially executed node (not needed as unreliable) • 0.9 — indicates data obtained from a finished node, but from a partially executed statement
<code>oids</code>	List of IDs of tables that were involved in the prediction for this node.

F.2.4.2.4. aqo_query_stat

The `aqo_query_stat` view shows statistics on query execution, by query class. aqo uses this data when `auto_tuning` is enabled for a particular query class.

Table F.4. aqo_query_stat View

Data	Description
<code>fs</code>	The identifier of the feature space.
<code>dbid</code>	The identifier of the database.

Data	Description
<code>execution_time_with_aqo</code>	Array of execution times for queries run with aqo enabled.
<code>execution_time_without_aqo</code>	Array of execution times for queries run with aqo disabled.
<code>planning_time_with_aqo</code>	Array of planning times for queries run with aqo enabled.
<code>planning_time_without_aqo</code>	Array of planning times for queries run with aqo disabled.
<code>cardinality_error_with_aqo</code>	Array of cardinality estimation errors in the selected query plans with aqo enabled.
<code>cardinality_error_without_aqo</code>	Array of cardinality estimation errors in the selected query plans with aqo disabled.
<code>executions_with_aqo</code>	Number of queries run with aqo enabled.
<code>executions_without_aqo</code>	Number of queries run with aqo disabled.

F.2.4.3. Functions

aqo adds several functions to Postgres Pro catalog.

F.2.4.3.1. Storage Management Functions

Important

Functions `aqo_queries_update`, `aqo_query_texts_update`, `aqo_query_stat_update`, `aqo_data_update` and `aqo_data_delete` modify data files underlying aqo views. Therefore, call these functions only if you understand the logic of adaptive query optimization.

`aqo_cleanup()` → setof integer

Removes data related to query classes that are linked (may be partially) with removed relations. Returns the number of removed feature spaces (classes) and feature subspaces. Insensitive to removing other objects.

`aqo_enable_class(fs bigint, dbid oid) → void`

Sets `learn_aqo`, `use_aqo` and `auto_tuning` (only in the intelligent mode) to true for the query class with the specified `fs` and `dbid`. You can set `dbid` to NULL instead of the ID of the current database.

`aqo_disable_class(fs bigint, dbid oid) → void`

Sets `learn_aqo`, `use_aqo` and `auto_tuning` to false for the query class with the specified `fs` and `dbid`. You can set `dbid` to NULL instead of the ID of the current database.

`aqo_drop_class(fs bigint, dbid oid) → integer`

Removes all data related to the specified query class and database from the aqo storage. You can set `dbid` to NULL instead of the ID of the current database. Returns the number of records removed from the aqo storage.

`aqo_reset(dbid oid) → bigint`

Removes records from the specified database: machine learning data, query texts, statistics and query class preferences. If `dbid` is omitted, removes the data from the current database. If `dbid` is NULL, removes all records from the aqo storage. Returns the number of records removed.

`aqo_queries_update` (*fs* bigint, *dbid* oid, *learn_aqo* boolean, *use_aqo* boolean, *auto_tuning* boolean) → boolean

Updates or inserts a record in a data file underlying the [aqo_queries](#) view for the specified *fs* and *dbid*. You can set *dbid* to NULL instead of the ID of the current database. NULL values for parameters being set mean leave them as is. Note that records with a zero value of *fs* or *dbid* cannot be updated. Returns *false* in case of error, *true* otherwise.

`aqo_query_texts_update` (*fs* bigint, *dbid* oid, *query_text* text) → boolean

Updates or inserts a record in a data file underlying the [aqo_query_texts](#) view for the specified *fs* and *dbid*. You can set *dbid* to NULL instead of the ID of the current database. Note that records with a zero value of *fs* or *dbid* cannot be updated. Returns *false* in case of error, *true* otherwise.

`aqo_query_stat_update` (*fs* bigint, *dbid* oid, *execution_time_with_aqo* double precision[], *execution_time_without_aqo* double precision[], *planning_time_with_aqo* double precision[], *planning_time_without_aqo* double precision[], *cardinality_error_with_aqo* double precision[], *cardinality_error_without_aqo* double precision[], *executions_with_aqo* bigint[], *executions_without_aqo* bigint[]) → boolean

Updates or inserts a record in a data file underlying the [aqo_query_stat](#) view for the specified *fs* and *dbid*. You cannot update a record for a common feature space, i.e., for *fs* or *dbid* equal to zero. You can set *dbid* to NULL instead of the ID of the current database. Returns *false* in case of error, *true* otherwise.

`aqo_data_update` (*fs* bigint, *fss* integer, *dbid* oid, *delta_rows* boolean, *nfeatures* integer, *features* double precision[][] , *targets* double precision[], *reliability* double precision[], *oids* oid[]) → boolean

Updates or inserts a record in a data file underlying the [aqo_data](#) view for the specified *fs*, *fss* and *dbid*. You can set *dbid* to NULL instead of the ID of the current database. If you set *delta_rows* to NULL, the value of `aqo.delta_rows` is used. Returns *false* in case of error, *true* otherwise.

`aqo_data_delete` (*fs* bigint, *fss* integer, *dbid* oid, *delta_rows* boolean) → boolean

Removes a record from a data file underlying the [aqo_data](#) view for the specified *fs*, *fss* and *dbid*. You can set *dbid* to NULL instead of the ID of the current database. If you omit *delta_rows* or set it to NULL, the value of `aqo.delta_rows` is used. Returns *false* in case of error, *true* otherwise.

F.2.4.3.2. Memory Management Functions

`aqo_memory_usage` () → setof record

Shows allocated and used sizes of aqo memory contexts and hash tables. Returns a table:

name

Short description of the memory context or hash table

allocated_size

Total size of the allocated memory

used_size

Size of the currently used memory

F.2.4.3.3. Functions for Analytics

`aqo_cardinality_error` (*controlled* boolean) → setof record

Shows the cardinality error for the last execution of queries. If *controlled* is true, shows queries executed with aqo enabled. If *controlled* is false, shows queries that were executed with aqo disabled, but that have collected aqo statistics. Returns a table:

num

Sequential number

dbid

The identifier of the database

fs

The identifier of the feature space. Can be zero or the basic hash.

error

aqo error calculated on query plan nodes

nexecs

Number of executions of queries associated with this *fs*

`aqo_execution_time (controlled boolean) → setof record`

Shows the execution time for queries. If *controlled* is true, shows the execution time of the last execution with aqo enabled. If *controlled* is false, returns the average execution time for all logged executions with aqo disabled. Execution time without aqo can be collected when `aqo.mode = intelligent` or `aqo.force_collect_stat = on`. Returns a table:

num

Sequential number

dbid

The identifier of the database

fs

The identifier of the feature space. Can be zero or the basic hash.

exec_time

If *controlled* = true, last query execution time with aqo, otherwise, average execution time for all executions without aqo

nexecs

Number of executions of queries associated with this *fs*

F.2.5. Examples

Example F.1. Learning on a Query (Basic Mode)

Consider optimization of a query using aqo.

When the query is executed for the first time, it is missing in tables underlying aqo views. So there is no data for predicting with aqo for each plan node, and “AQO not used” lines appear in the `EXPLAIN` output:

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT bp.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';
```

QUERY PLAN

```
-----
Hash Join  (cost=63796.33..140286.73 rows=477728 width=25) (actual rows=2215 loops=1)
  AQO not used, fss=8598194613120045129
```

Additional Supplied Modules and Extensions Shipped in

postgrespro-sdm-17-contrib

```

Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
-> Seq Scan on boarding_passes bp (cost=0.00..32894.95 rows=1894295 width=25)
(actual rows=1894295 loops=1)
    AJO not used, fss=1362775811343989307
-> Hash (cost=51379.43..51379.43 rows=595260 width=22) (actual rows=468255
loops=1)
    Buckets: 131072 Batches: 8 Memory Usage: 4223kB
-> Hash Join (cost=1859.80..51379.43 rows=595260 width=22) (actual
rows=468255 loops=1)
    AJO not used, fss=-7651474063207585780
    Hash Cond: (tf.flight_id = f.flight_id)
-> Seq Scan on ticket_flights tf (cost=0.00..43323.35 rows=2360335
width=18) (actual rows=2360335 loops=1)
    AJO not used, fss=-6410966714754547713
-> Hash (cost=1652.80..1652.80 rows=16560 width=4) (actual rows=16340
loops=1)
    Buckets: 32768 Batches: 1 Memory Usage: 831kB
-> Seq Scan on flights f (cost=0.00..1652.80 rows=16560 width=4)
(actual rows=16340 loops=1)
    AJO not used, fss=-1289471166524579716
    Filter: (scheduled_departure > '2017-08-15
22:00:00+07':::timestamp with time zone)
    Rows Removed by Filter: 49324

Using ajo: true
AJO mode: LEARN
AJO advanced: OFF
Query hash: 402936317627943257
JOINS: 2
(23 rows)

```

If there is no information on a certain node in the [ajo_data](#) view, ajo will add the appropriate record there for future learning and predictions except for nodes with `fss=0` in the EXPLAIN output. As each of features and targets in the `ajo_data` view is a logarithm to base *e*, to get the actual value, raise *e* to this power. For example: `exp(7.703007682479236)`:

```

demo=# select * from ajo_data;
 fs |          fss          | dbid | delta_rows | nfeatures |          features
----+-----+-----+-----+-----+-----+-----
 0 | 1362775811343989307 | 16429 | f          | 0         | {14.454357295615447} | {1} | {16452}
 0 | -6410966714754547713 | 16429 | f          | 0         | {14.674314116080508} | {1} | {16479}
 0 | -1289471166524579716 | 16429 | f          | 1         | {9.701371368413994}  | {1} | {16458}
 0 | 8598194613120045129 | 16429 | f          | 2         | {-13.492416828684513,-1.3775704575284085} | {7.703007682479236} | {1} | {16479,16458,16452}
 0 | -7651474063207585780 | 16429 | f          | 2         | {-11.092306109090387,-1.3775704575284085} | {13.056768298305919} | {1} | {16479,16458}
(5 rows)

```

When the query is executed for the second time, ajo recognizes the query and makes a prediction. Pay attention to the cardinality predicted by ajo and the value of ajo error ("error=0%").

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
```

**Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib**

```
SELECT bp.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';
```

QUERY PLAN

```
Hash Join  (cost=61144.50..136890.86 rows=2215 width=25) (actual rows=2215 loops=1)
  AQO: rows=2215, error=0%, fss=8598194613120045129
  Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
    -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295 width=25)
        (actual rows=1894295 loops=1)
        AQO: rows=1894295, error=0%, fss=1362775811343989307
    -> Hash  (cost=51376.68..51376.68 rows=468255 width=22) (actual rows=468255
        loops=1)
        Buckets: 131072  Batches: 4  Memory Usage: 7438kB
        -> Hash Join  (cost=1857.05..51376.68 rows=468255 width=22) (actual
            rows=468255 loops=1)
            AQO: rows=468255, error=0%, fss=-7651474063207585780
            Hash Cond: (tf.flight_id = f.flight_id)
            -> Seq Scan on ticket_flights tf  (cost=0.00..43323.35 rows=2360335
                width=18) (actual rows=2360335 loops=1)
                AQO: rows=2360335, error=0%, fss=-6410966714754547713
            -> Hash  (cost=1652.80..1652.80 rows=16340 width=4) (actual rows=16340
                loops=1)
                Buckets: 16384  Batches: 1  Memory Usage: 703kB
                -> Seq Scan on flights f  (cost=0.00..1652.80 rows=16340 width=4)
                    (actual rows=16340 loops=1)
                    AQO: rows=16340, error=0%, fss=-1289471166524579716
                    Filter: (scheduled_departure > '2017-08-15
22:00:00+07'::timestamp with time zone)
                    Rows Removed by Filter: 49324

Using aqo: true
AQO mode: LEARN
AQO advanced: OFF
Query hash: 402936317627943257
JOINS: 2
(23 rows)
```

Let's change a constant in the query, and you will notice that the prediction is made with an error:

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT bp.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-10 15:00:00+00';
```

QUERY PLAN

```
Hash Join  (cost=61144.50..136890.86 rows=2215 width=25) (actual rows=111397 loops=1)
  AQO: rows=2215, error=-4929%, fss=8598194613120045129
  Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
    -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295 width=25)
        (actual rows=1894295 loops=1)
        AQO: rows=1894295, error=0%, fss=1362775811343989307
    -> Hash  (cost=51376.68..51376.68 rows=468255 width=22) (actual rows=577437
        loops=1)
        Buckets: 131072 (originally 131072)  Batches: 8 (originally 4)  Memory Usage:
7169kB
```

postgrespro-sdm-17-contrib

```
Using aqo: true
AQO mode: LEARN
AQO advanced: OFF
Query hash: 402936317627943257
JOINS: 2
(23 rows)
```

However, instead of recalculating `features` and `targets`, `aqo` added new values of selectivity and cardinality for this query to `aqo_data`:

```

demo=# select * from aqo_data;
 fs |          fss          | dbid  | delta_rows | nfeatures |          |
          features          |      | oids      |           |          | targets
          | reliability |
-----+-----+-----+-----+-----+-----
0 | 1362775811343989307 | 16429 | f          |           | 0 | {14.454357295615447}
          | {1}          | {16452}
0 | -6410966714754547713 | 16429 | f          |           | 0 | {14.6743141116080508}
          | {1}          | {16479}
0 | -1289471166524579716 | 16429 | f          |           | 1 | {{-1.3775704575284085},
{-1.2287385463463019}}
{9.701371368413994,9.854297308345357} | {1,1} | {16458}
0 | 8598194613120045129 | 16429 | f          |           | 2 |
{{-13.492416828684513,-1.3775704575284085},{-13.492416828684513,-1.2287385463463019}}
| {7.703007682479236,11.620855676130656} | {1,1} | {16479,1645
8,16452}
0 | -7651474063207585780 | 16429 | f          |           | 2 |
{{-11.092306109090387,-1.3775704575284085},{-11.092306109090387,-1.2287385463463019}}
| {13.056768298305919,13.266354624518149} | {1,1} | {16479,1645
8}
(5 rows)

```

Now the prediction has a small error of about 3%, which can be explained by a calculation error:

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT bp.*
```

**Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib**

```
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-10 15:00:00+00';
```

QUERY PLAN

```
-----
Hash Join  (cost=63438.72..139820.11 rows=108541 width=25) (actual rows=111397
loops=1)
  AQP: rows=108541, error=-3%, fss=8598194613120045129
  Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
    -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295 width=25)
    (actual rows=1894295 loops=1)
      AQP: rows=1894295, error=0%, fss=1362775811343989307
    -> Hash  (cost=51410.19..51410.19 rows=576635 width=22) (actual rows=577437
loops=1)
      Buckets: 131072  Batches: 8  Memory Usage: 4966kB
      -> Hash Join  (cost=1890.56..51410.19 rows=576635 width=22) (actual
rows=577437 loops=1)
        AQP: rows=576635, error=-0%, fss=-7651474063207585780
        Hash Cond: (tf.flight_id = f.flight_id)
          -> Seq Scan on ticket_flights tf  (cost=0.00..43323.35 rows=2360335
width=18) (actual rows=2360335 loops=1)
            AQP: rows=2360335, error=0%, fss=-6410966714754547713
          -> Hash  (cost=1652.80..1652.80 rows=19021 width=4) (actual rows=19040
loops=1)
            Buckets: 32768  Batches: 1  Memory Usage: 926kB
            -> Seq Scan on flights f  (cost=0.00..1652.80 rows=19021 width=4)
            (actual rows=19040 loops=1)
              AQP: rows=19021, error=-0%, fss=-1289471166524579716
              Filter: (scheduled_departure > '2017-08-10
22:00:00+07'::timestamp with time zone)
              Rows Removed by Filter: 46624

Using aqp: true
AQP mode: LEARN
AQP advanced: OFF
Query hash: 402936317627943257
JOINS: 2
(23 rows)
```

We can modify the query by adding some table to the JOIN list. In this case, aqp will predict the cardinality of nodes on which it learned before.

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT t.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
JOIN tickets t ON t.ticket_no = tf.ticket_no
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';
```

QUERY PLAN

```
-----
Nested Loop  (cost=61158.91..134296.78 rows=2273 width=104) (actual rows=2215 loops=1)
  AQP not used, fss=-8581941154270057231
  -> Hash Join  (cost=61158.48..133208.83 rows=2273 width=28) (actual rows=2215
loops=1)
    AQP: rows=2273, error=3%, fss=8598194613120045129
    Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
      -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295
width=18) (actual rows=1894295 loops=1)
```


Additional Supplied Modules and Extensions Shipped in

postgrespro-sdm-17-contrib

```
AQO: rows=1894295, error=0%, fss=1362775811343989307
-> Hash (cost=51376.89..51376.89 rows=468906 width=22) (actual rows=468255
loops=1)
    Buckets: 131072  Batches: 4  Memory Usage: 7438kB
    -> Hash Join (cost=1857.26..51376.89 rows=468906 width=22) (actual
rows=468255 loops=1)
        AQO: rows=468906, error=0%, fss=-7651474063207585780
        Hash Cond: (tf.flight_id = f.flight_id)
        -> Seq Scan on ticket_flights tf (cost=0.00..43323.35
rows=2360335 width=18) (actual rows=2360335 loops=1)
            AQO: rows=2360335, error=0%, fss=-6410966714754547713
            -> Hash (cost=1652.80..1652.80 rows=16357 width=4) (actual
rows=16340 loops=1)
                Buckets: 16384  Batches: 1  Memory Usage: 703kB
                -> Seq Scan on flights f (cost=0.00..1652.80 rows=16357
width=4) (actual rows=16340 loops=1)
                    AQO: rows=16357, error=0%, fss=-1289471166524579716
                    Filter: (scheduled_departure > '2017-08-15
22:00:00+07':::timestamp with time zone)
                        Rows Removed by Filter: 49324
                    -> Index Scan using tickets_pkey on tickets t (cost=0.42..0.48 rows=1 width=104)
(actual rows=1 loops=2215)
                        AQO not used, fss=2731022528523952664
                        Index Cond: (ticket_no = bp.ticket_no)
Using aqo: true
AQO mode: LEARN
AQO advanced: OFF
Query hash: 7809046947949890015
JOINS: 3
(28 rows)
```

Example F.2. Using the [aqo_query_stat](#) View

The `aqo_query_stats` view shows statistics on the query planning time, query execution time and cardinality error. Based on this data you can make a decision whether to use aqo predictions for different query classes.

Let's query the `aqo_query_stats` view:

```
demo=# SELECT * FROM aqo_query_stat \gx
-[ RECORD 1 ]-----+
fs                  | 7809046947949890015
dbid                | 16429
execution_time_with_aqo | {1.039218233,0.925258453,0.831166925,0.779602353}
execution_time_without_aqo | {1.022052611,0.936486619}
planning_time_with_aqo | {0.003305339,0.002129048,0.002538877,0.002142972}
planning_time_without_aqo | {0.000767553,0.000711208}
cardinality_error_with_aqo |
{0.4854215265638894,0,1.1711726076352047,0.007732205169478082}
cardinality_error_without_aqo | {0.4854215265638894,1.571562511977072}
executions_with_aqo | 4
executions_without_aqo | 2
```

The retrieved data is for the query from [Example F.1](#), which was executed once without aqo for each of the parameters `f.scheduled_departure > '2017-08-10 15:00:00+00'` and `f.scheduled_departure > '2017-08-15 15:00:00+00'` and twice with aqo for each of these parameters. It is clear that with aqo, the cardinality error decreases to 0.0077, while the minimum cardinality error without aqo is 0.4854. Besides, the execution time with aqo is lower than without it. So the conclusion is that aqo learns well on this query, and the prediction can be used for this query class.

Example F.3. Using aqo in the Advanced Mode

The advanced mode allows a more flexible control over aqo. When this mode is activated, that is,

```
demo=# SET aqo.advanced = on;
```

aqo will collect the machine learning data separately for each query executed. For example:

```
demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT bp.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';
                                QUERY PLAN
-----
Hash Join  (cost=63796.33..140286.73 rows=477728 width=25) (actual rows=2215 loops=1)
  AQO not used, fss=8598194613120045129
  Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
    -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295 width=25)
        (actual rows=1894295 loops=1)
        AQO not used, fss=1362775811343989307
        -> Hash  (cost=51379.43..51379.43 rows=595260 width=22) (actual rows=468255
            loops=1)
            Buckets: 131072  Batches: 8  Memory Usage: 4223kB
            -> Hash Join  (cost=1859.80..51379.43 rows=595260 width=22) (actual
                rows=468255 loops=1)
                AQO not used, fss=-7651474063207585780
                Hash Cond: (tf.flight_id = f.flight_id)
                -> Seq Scan on ticket_flights tf  (cost=0.00..43323.35 rows=2360335
                    width=18) (actual rows=2360335 loops=1)
                    AQO not used, fss=-6410966714754547713
                    -> Hash  (cost=1652.80..1652.80 rows=16560 width=4) (actual rows=16340
                        loops=1)
                        Buckets: 32768  Batches: 1  Memory Usage: 831kB
                        -> Seq Scan on flights f  (cost=0.00..1652.80 rows=16560 width=4)
                            (actual rows=16340 loops=1)
                            AQO not used, fss=-1289471166524579716
                            Filter: (scheduled_departure > '2017-08-15
                                22:00:00+07'::timestamp with time zone)
                                Rows Removed by Filter: 49324

Using aqo: true
AQO mode: LEARN
AQO advanced: ON
Query hash: 402936317627943257
JOINS: 2
(23 rows)
```

Now this query is stored in [aqo_data](#) with a non-zero fs (fs is equal to the query hash by default):

```
demo=# SELECT * FROM aqo_data;
      fs      |      fss      | dbid | delta_rows | nfeatures |      |      oids
      features |      targets  |      | reliability |      |
-----+-----+-----+-----+-----+-----+-----
402936317627943257 | -6410966714754547713 | 16429 | f          |          | 0 | {16479}
                  | {14.674314116080508} | {1}   |          |          | 0 | {16452}
402936317627943257 | 1362775811343989307 | 16429 | f          |          | 0 | {16452}
                  | {14.454357295615447} | {1}   |          |          | 0 | {16452}
```

**Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib**

```

402936317627943257 | -1289471166524579716 | 16429 | f | 1 |
{{-1.3775704575284085}} | {9.701371368413994} | {1} |
{16458}
402936317627943257 | 8598194613120045129 | 16429 | f | 2 |
{{-13.492416828684513,-1.3775704575284085}} | {7.703007682479236} | {1} |
{16479,16458,16452}
402936317627943257 | -7651474063207585780 | 16429 | f | 2 |
{{-11.092306109090387,-1.3775704575284085}} | {13.056768298305919} | {1} |
{16479,16458}
(5 rows)

```

We can make a few settings individually for this query. These are values of `learn_aqo`, `use_aqo` and `auto_tuning` in the [aqo_queries](#) view:

```

demo=# SELECT * FROM aqo_queries;
      fs      | dbid | learn_aqo | use_aqo | auto_tuning | smart_timeout |
count_increase_timeout
-----+-----+-----+-----+-----+-----+-----
+-----+
402936317627943257 | 16429 | t          | t          | f          |              | 0 |
0
0 | 0 | f          | f          | f          |              | 0 |
0
(2 rows)

```

Let's set `use_aqo` to false:

```

demo=# SELECT aqo_queries_update(402936317627943257, NULL, NULL, false, NULL);
aqo_queries_update
-----
t
(1 row)

```

Now we change a constant in the query:

```

demo=# EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT bp.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
WHERE f.scheduled_departure > '2017-08-10 15:00:00+00';

```

QUERY PLAN

```

-----
Hash Join  (cost=65822.16..142872.60 rows=554378 width=25) (actual rows=111397
loops=1)
  AQO not used, fss=0
  Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))
    -> Seq Scan on boarding_passes bp  (cost=0.00..32894.95 rows=1894295 width=25)
(actual rows=1894295 loops=1)
      AQO not used, fss=0
      -> Hash  (cost=51412.64..51412.64 rows=690768 width=22) (actual rows=577437
loops=1)
        Buckets: 131072  Batches: 8  Memory Usage: 4966kB
        -> Hash Join  (cost=1893.01..51412.64 rows=690768 width=22) (actual
rows=577437 loops=1)
          AQO not used, fss=0
          Hash Cond: (tf.flight_id = f.flight_id)
            -> Seq Scan on ticket_flights tf  (cost=0.00..43323.35 rows=2360335
width=18) (actual rows=2360335 loops=1)
              AQO not used, fss=0

```

**Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib**

```

-> Hash  (cost=1652.80..1652.80 rows=19217 width=4) (actual rows=19040
loops=1)
      Buckets: 32768  Batches: 1  Memory Usage: 926kB
      -> Seq Scan on flights f  (cost=0.00..1652.80 rows=19217 width=4)
(actual rows=19040 loops=1)
          AQP not used, fss=0
          Filter: (scheduled_departure > '2017-08-10
22:00:00+07':::timestamp with time zone)
          Rows Removed by Filter: 46624

Using aqo: false
AQO mode: LEARN
AQO advanced: ON
Query hash: 402936317627943257
JOINS: 2
(23 rows)

```

aqo was not used for this query, but there is new data in the aqo_data view:

```

demo=# SELECT * FROM aqo_data;
      fs          |          fss          | dbid | delta_rows | nfeatures |
      features          |
      targets          | reliability |      oids
-----+-----+-----+-----+-----+
402936317627943257 | -6410966714754547713 | 16429 | f          |          0 |
{14.674314116080508} | {1}          | {16479}
402936317627943257 | 1362775811343989307 | 16429 | f          |          0 |
{14.454357295615447} | {1}          | {16452}
402936317627943257 | -1289471166524579716 | 16429 | f          |          1 |
{{-1.3775704575284085},{-1.2287385463463019}}
| {9.701371368413994,9.854297308345357} | {1,1}          | {16458}
402936317627943257 | 8598194613120045129 | 16429 | f          |          2 |
{{-13.492416828684513,-1.3775704575284085},{-13.492416828684513,-1.2287385463463019}}
| {7.703007682479236,11.620855676130656} | {1,1}          | {16479,16458,16452}
402936317627943257 | -7651474063207585780 | 16429 | f          |          2 |
{{-11.092306109090387,-1.3775704575284085},{-11.092306109090387,-1.2287385463463019}}
| {13.056768298305919,13.266354624518149} | {1,1}          | {16479,16458}
(5 rows)

```

The use_aqo setting does not apply to other queries. After executing another query twice, we can see that aqo learns on it and makes prediction for it:

```

EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT t.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
JOIN tickets t ON t.ticket_no = tf.ticket_no
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';

```

QUERY PLAN

```

Nested Loop  (cost=61144.93..134251.05 rows=2215 width=104) (actual rows=2215 loops=1)
  AQP: rows=2215, error=0%, fss=-8581941154270057231
  -> Hash Join  (cost=61144.50..133190.86 rows=2215 width=28) (actual rows=2215
loops=1)
      AQP: rows=2215, error=0%, fss=8598194613120045129
      Hash Cond: ((bp.flight_id = f.flight_id) AND (bp.ticket_no = tf.ticket_no))

```

Additional Supplied Modules and Extensions Shipped in

postgrespro-sdm-17-contrib

```
-> Seq Scan on boarding_passes bp (cost=0.00..32894.95 rows=1894295
width=18) (actual rows=1894295 loops=1)
    AQP: rows=1894295, error=0%, fss=1362775811343989307
-> Hash (cost=51376.68..51376.68 rows=468255 width=22) (actual rows=468255
loops=1)
    Buckets: 131072 Batches: 4 Memory Usage: 7438kB
-> Hash Join (cost=1857.05..51376.68 rows=468255 width=22) (actual
rows=468255 loops=1)
    AQP: rows=468255, error=0%, fss=-7651474063207585780
    Hash Cond: (tf.flight_id = f.flight_id)
-> Seq Scan on ticket_flights tf (cost=0.00..43323.35
rows=2360335 width=18) (actual rows=2360335 loops=1)
    AQP: rows=2360335, error=0%, fss=-6410966714754547713
-> Hash (cost=1652.80..1652.80 rows=16340 width=4) (actual
rows=16340 loops=1)
    Buckets: 16384 Batches: 1 Memory Usage: 703kB
-> Seq Scan on flights f (cost=0.00..1652.80 rows=16340
width=4) (actual rows=16340 loops=1)
    AQP: rows=16340, error=0%, fss=-1289471166524579716
    Filter: (scheduled_departure > '2017-08-15
22:00:00+07':::timestamp with time zone)
    Rows Removed by Filter: 49324
-> Index Scan using tickets_pkey on tickets t (cost=0.42..0.48 rows=1 width=104)
(actual rows=1 loops=2215)
    AQP not used, fss=2731022528523952664
    Index Cond: (ticket_no = bp.ticket_no)
Using aqp: true
AQP mode: LEARN
AQP advanced: ON
Query hash: 7809046947949890015
JOINS: 3
(28 rows)
```

Example F.4. Using the Sandbox Mode

```
SET aqp.sandbox = ON;
SET aqp.enable = ON;
SET aqp.advanced = OFF;
-- Clean up the sandbox knowledge base without touching the main data
SELECT aqp_reset();

EXPLAIN (ANALYZE, SUMMARY OFF, TIMING OFF)
SELECT t.*
FROM ticket_flights tf
JOIN flights f ON f.flight_id = tf.flight_id
JOIN boarding_passes bp ON bp.ticket_no = tf.ticket_no AND bp.flight_id = tf.flight_id
JOIN tickets t ON t.ticket_no = tf.ticket_no
WHERE f.scheduled_departure > '2017-08-15 15:00:00+00';

-- Be executing the previous query until plans get stabilized
...

-- Copy data obtained with aqp.advanced = OFF from sandbox
CREATE TABLE aqp_data_sandbox AS SELECT * FROM aqp_data;
SET aqp.sandbox = OFF;
SELECT aqp_data_update (fs, fss, dbid, delta_rows, nfeatures, features, targets,
    reliability, oids)
FROM aqp_data_sandbox WHERE fs = 0;
DROP TABLE aqp_data_sandbox;
```

F.2.6. Author

Oleg Ivanov

F.3. auto_explain — log execution plans of slow queries

The `auto_explain` module provides a means for logging execution plans of slow statements automatically, without having to run `EXPLAIN` by hand. This is especially helpful for tracking down un-optimized queries in large applications.

The module provides no SQL-accessible functions. To use it, simply load it into the server. You can load it into an individual session:

```
LOAD 'auto_explain';
```

(You must be superuser to do that.) More typical usage is to preload it into some or all sessions by including `auto_explain` in [session_preload_libraries](#) or [shared_preload_libraries](#) in `postgresql.conf`. Then you can track unexpectedly slow queries no matter when they happen. Of course there is a price in overhead for that.

F.3.1. Configuration Parameters

There are several configuration parameters that control the behavior of `auto_explain`. Note that the default behavior is to do nothing, so you must set at least `auto_explain.log_min_duration` if you want any results.

`auto_explain.log_min_duration (integer)`

`auto_explain.log_min_duration` is the minimum statement execution time, in milliseconds, that will cause the statement's plan to be logged. Setting this to 0 logs all plans. -1 (the default) disables logging of plans. For example, if you set it to 250ms then all statements that run 250ms or longer will be logged. Only superusers can change this setting.

`auto_explain.log_parameter_max_length (integer)`

`auto_explain.log_parameter_max_length` controls the logging of query parameter values. A value of -1 (the default) logs the parameter values in full. 0 disables logging of parameter values. A value greater than zero truncates each parameter value to that many bytes. Only superusers can change this setting.

`auto_explain.log_analyze (boolean)`

`auto_explain.log_analyze` causes `EXPLAIN ANALYZE` output, rather than just `EXPLAIN` output, to be printed when an execution plan is logged. This parameter is off by default. Only superusers can change this setting.

Note

When this parameter is on, per-plan-node timing occurs for all statements executed, whether or not they run long enough to actually get logged. This can have an extremely negative impact on performance. Turning off `auto_explain.log_timing` ameliorates the performance cost, at the price of obtaining less information.

`auto_explain.log_buffers (boolean)`

`auto_explain.log_buffers` controls whether buffer usage statistics are printed when an execution plan is logged; it's equivalent to the `BUFFERS` option of `EXPLAIN`. This parameter has no effect unless `auto_explain.log_analyze` is enabled. This parameter is off by default. Only superusers can change this setting.

`auto_explain.log_wal (boolean)`

`auto_explain.log_wal` controls whether WAL usage statistics are printed when an execution plan is logged; it's equivalent to the `WAL` option of `EXPLAIN`. This parameter has no effect unless `auto_explain.log_analyze` is enabled. This parameter is off by default. Only superusers can change this setting.

`auto_explain.log_timing` (boolean)

`auto_explain.log_timing` controls whether per-node timing information is printed when an execution plan is logged; it's equivalent to the `TIMING` option of `EXPLAIN`. The overhead of repeatedly reading the system clock can slow down queries significantly on some systems, so it may be useful to set this parameter to off when only actual row counts, and not exact times, are needed. This parameter has no effect unless `auto_explain.log_analyze` is enabled. This parameter is on by default. Only superusers can change this setting.

`auto_explain.log_triggers` (boolean)

`auto_explain.log_triggers` causes trigger execution statistics to be included when an execution plan is logged. This parameter has no effect unless `auto_explain.log_analyze` is enabled. This parameter is off by default. Only superusers can change this setting.

`auto_explain.log_verbose` (boolean)

`auto_explain.log_verbose` controls whether verbose details are printed when an execution plan is logged; it's equivalent to the `VERBOSE` option of `EXPLAIN`. This parameter is off by default. Only superusers can change this setting.

`auto_explain.log_settings` (boolean)

`auto_explain.log_settings` controls whether information about modified configuration options is printed when an execution plan is logged. Only options affecting query planning with value different from the built-in default value are included in the output. This parameter is off by default. Only superusers can change this setting.

`auto_explain.log_format` (enum)

`auto_explain.log_format` selects the `EXPLAIN` output format to be used. The allowed values are `text`, `xml`, `json`, and `yaml`. The default is `text`. Only superusers can change this setting.

`auto_explain.log_level` (enum)

`auto_explain.log_level` selects the log level at which `auto_explain` will log the query plan. Valid values are `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `INFO`, `NOTICE`, `WARNING`, and `LOG`. The default is `LOG`. Only superusers can change this setting.

`auto_explain.log_nested_statements` (boolean)

`auto_explain.log_nested_statements` causes nested statements (statements executed inside a function) to be considered for logging. When it is off, only top-level query plans are logged. This parameter is off by default. Only superusers can change this setting.

`auto_explain.sample_rate` (real)

`auto_explain.sample_rate` causes `auto_explain` to only explain a fraction of the statements in each session. The default is 1, meaning explain all the queries. In case of nested statements, either all will be explained or none. Only superusers can change this setting.

In ordinary usage, these parameters are set in `postgresql.conf`, although superusers can alter them on-the-fly within their own sessions. Typical usage might be:

```
# postgresql.conf
session_preload_libraries = 'auto_explain'

auto_explain.log_min_duration = '3s'
```

F.3.2. Example

```
postgres=# LOAD 'auto_explain';
postgres=# SET auto_explain.log_min_duration = 0;
```



```
postgres=# SET auto_explain.log_analyze = true;
postgres=# SELECT count(*)
           FROM pg_class, pg_index
           WHERE oid = indrelid AND indisunique;
```

This might produce log output such as:

```
LOG:  duration: 0.196 ms planning: 0.548 ms plan:
       Query Text: SELECT count(*)
                   FROM pg_class, pg_index
                   WHERE oid = indrelid AND indisunique;
       Aggregate  (cost=16.79..16.80 rows=1 width=0) (actual time=3.626..3.627 rows=1
loops=1)
       -> Hash Join  (cost=4.17..16.55 rows=92 width=0) (actual time=3.349..3.594 rows=92
loops=1)
           Hash Cond: (pg_class.oid = pg_index.indrelid)
           -> Seq Scan on pg_class  (cost=0.00..9.55 rows=255 width=4) (actual
time=0.016..0.140 rows=255 loops=1)
           -> Hash  (cost=3.02..3.02 rows=92 width=4) (actual time=3.238..3.238 rows=92
loops=1)
               Buckets: 1024  Batches: 1  Memory Usage: 4kB
               -> Seq Scan on pg_index  (cost=0.00..3.02 rows=92 width=4) (actual
time=0.008..3.187 rows=92 loops=1)
                   Filter: indisunique
```

Note that Postgres Pro added a new field to the output to show planning time.

F.3.3. Authors

Takahiro Itagaki <itagaki.takahiro@oss.ntt.co.jp>

Postgres Professional, Moscow, Russia

F.4. dblink — connect to other Postgres Pro databases

`dblink` is a module that supports connections to other Postgres Pro databases from within a database session.

`dblink` can report the following wait events under the wait event type `Extension`.

`DblinkConnect`

Waiting to establish a connection to a remote server.

`DblinkGetConnect`

Waiting to establish a connection to a remote server when it could not be found in the list of already-opened connections.

`DblinkGetResult`

Waiting to receive the results of a query from a remote server.

See also [postgres_fdw](#), which provides roughly the same functionality using a more modern and standards-compliant infrastructure.

dblink_connect

`dblink_connect` — opens a persistent connection to a remote database

Synopsis

```
dblink_connect(text connstr) returns text
dblink_connect(text connname, text connstr) returns text
```

Description

`dblink_connect()` establishes a connection to a remote Postgres Pro database. The server and database to be contacted are identified through a standard libpq connection string. Optionally, a name can be assigned to the connection. Multiple named connections can be open at once, but only one unnamed connection is permitted at a time. The connection will persist until closed or until the database session is ended.

The connection string may also be the name of an existing foreign server. It is recommended to use the foreign-data wrapper `dblink_fdw` when defining the foreign server. See the example below, as well as [CREATE SERVER](#).

Arguments

connname

The name to use for this connection; if omitted, an unnamed connection is opened, replacing any existing unnamed connection.

connstr

libpq-style connection info string, for example `hostaddr=127.0.0.1 port=5432 dbname=mydb user=postgres password=mypasswd options=-csearch_path=`. For details see [Section 35.1.1](#). Alternatively, the name of a foreign server.

Return Value

Returns status, which is always OK (since any error causes the function to throw an error instead of returning).

Notes

If untrusted users have access to a database that has not adopted a [secure schema usage pattern](#), begin each session by removing publicly-writable schemas from `search_path`. One could, for example, add `options=-csearch_path=` to *connstr*. This consideration is not specific to `dblink`; it applies to every interface for executing arbitrary SQL commands.

Only superusers may use `dblink_connect` to create non-password-authenticated and non-GSSAPI-authenticated connections. If non-superusers need this capability, use `dblink_connect_u` instead.

It is unwise to choose connection names that contain equal signs, as this opens a risk of confusion with connection info strings in other `dblink` functions.

Examples

```
SELECT dblink_connect('dbname=postgres options=-csearch_path=');
 dblink_connect
-----
      OK
(1 row)

SELECT dblink_connect('myconn', 'dbname=postgres options=-csearch_path=');
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

```
dblink_connect
-----
OK
(1 row)

-- FOREIGN DATA WRAPPER functionality
-- Note: local connection must require password authentication for this to work
properly
--      Otherwise, you will receive the following error from dblink_connect():
--      ERROR:  password is required
--      DETAIL:  Non-superuser cannot connect if the server does not request a
password.
--      HINT:   Target server's authentication method must be changed.

CREATE SERVER fdtest FOREIGN DATA WRAPPER dblink_fdw OPTIONS (hostaddr '127.0.0.1',
    dbname 'contrib_regression');

CREATE USER regress_dblink_user WITH PASSWORD 'secret';
CREATE USER MAPPING FOR regress_dblink_user SERVER fdtest OPTIONS (user
    'regress_dblink_user', password 'secret');
GRANT USAGE ON FOREIGN SERVER fdtest TO regress_dblink_user;
GRANT SELECT ON TABLE foo TO regress_dblink_user;

\set ORIGINAL_USER :USER
\c - regress_dblink_user
SELECT dblink_connect('myconn', 'fdtest');
dblink_connect
-----
OK
(1 row)

SELECT * FROM dblink('myconn', 'SELECT * FROM foo') AS t(a int, b text, c text[]);
 a | b |          c
----+---+-----
 0 | a | {a0,b0,c0}
 1 | b | {a1,b1,c1}
 2 | c | {a2,b2,c2}
 3 | d | {a3,b3,c3}
 4 | e | {a4,b4,c4}
 5 | f | {a5,b5,c5}
 6 | g | {a6,b6,c6}
 7 | h | {a7,b7,c7}
 8 | i | {a8,b8,c8}
 9 | j | {a9,b9,c9}
10 | k | {a10,b10,c10}
(11 rows)

\c - :ORIGINAL_USER
REVOKE USAGE ON FOREIGN SERVER fdtest FROM regress_dblink_user;
REVOKE SELECT ON TABLE foo FROM regress_dblink_user;
DROP USER MAPPING FOR regress_dblink_user SERVER fdtest;
DROP USER regress_dblink_user;
DROP SERVER fdtest;
```

dblink_connect_u

`dblink_connect_u` — opens a persistent connection to a remote database, insecurely

Synopsis

```
dblink_connect_u(text connstr) returns text  
dblink_connect_u(text connname, text connstr) returns text
```

Description

`dblink_connect_u()` is identical to `dblink_connect()`, except that it will allow non-superusers to connect using any authentication method.

If the remote server selects an authentication method that does not involve a password, then impersonation and subsequent escalation of privileges can occur, because the session will appear to have originated from the user as which the local Postgres Pro server runs. Also, even if the remote server does demand a password, it is possible for the password to be supplied from the server environment, such as a `~/.pgpass` file belonging to the server's user. This opens not only a risk of impersonation, but the possibility of exposing a password to an untrustworthy remote server. Therefore, `dblink_connect_u()` is initially installed with all privileges revoked from `PUBLIC`, making it un-callable except by superusers. In some situations it may be appropriate to grant `EXECUTE` permission for `dblink_connect_u()` to specific users who are considered trustworthy, but this should be done with care. It is also recommended that any `~/.pgpass` file belonging to the server's user *not* contain any records specifying a wildcard host name.

For further details see `dblink_connect()`.

dblink_disconnect

`dblink_disconnect` — closes a persistent connection to a remote database

Synopsis

```
dblink_disconnect() returns text  
dblink_disconnect(text connname) returns text
```

Description

`dblink_disconnect()` closes a connection previously opened by `dblink_connect()`. The form with no arguments closes an unnamed connection.

Arguments

connname

The name of a named connection to be closed.

Return Value

Returns status, which is always `OK` (since any error causes the function to throw an error instead of returning).

Examples

```
SELECT dblink_disconnect();  
 dblink_disconnect  
-----  
OK  
(1 row)  
  
SELECT dblink_disconnect('myconn');  
 dblink_disconnect  
-----  
OK  
(1 row)
```

dblink

`dblink` — executes a query in a remote database

Synopsis

```
dblink(text connname, text sql [, bool fail_on_error]) returns setof record
dblink(text connstr, text sql [, bool fail_on_error]) returns setof record
dblink(text sql [, bool fail_on_error]) returns setof record
```

Description

`dblink` executes a query (usually a `SELECT`, but it can be any SQL statement that returns rows) in a remote database.

When two `text` arguments are given, the first one is first looked up as a persistent connection's name; if found, the command is executed on that connection. If not found, the first argument is treated as a connection info string as for `dblink_connect`, and the indicated connection is made just for the duration of this command.

Arguments

connname

Name of the connection to use; omit this parameter to use the unnamed connection.

connstr

A connection info string, as previously described for `dblink_connect`.

sql

The SQL query that you wish to execute in the remote database, for example `select * from foo`.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a `NOTICE`, and the function returns no rows.

Return Value

The function returns the row(s) produced by the query. Since `dblink` can be used with any query, it is declared to return `record`, rather than specifying any particular set of columns. This means that you must specify the expected set of columns in the calling query — otherwise Postgres Pro would not know what to expect. Here is an example:

```
SELECT *
FROM dblink('dbname=mydb options=-csearch_path=',
            'select pronomame, prosrc from pg_proc')
AS t1(pronomame name, prosrc text)
WHERE pronomame LIKE 'bytea%';
```

The “alias” part of the `FROM` clause must specify the column names and types that the function will return. (Specifying column names in an alias is actually standard SQL syntax, but specifying column types is a Postgres Pro extension.) This allows the system to understand what `*` should expand to, and what `pronomame` in the `WHERE` clause refers to, in advance of trying to execute the function. At run time, an error will be thrown if the actual query result from the remote database does not have the same number of columns shown in the `FROM` clause. The column names need not match, however, and `dblink` does not insist on exact type matches either. It will succeed so long as the returned data strings are valid input for the column type declared in the `FROM` clause.

Notes

A convenient way to use `dblink` with predetermined queries is to create a view. This allows the column type information to be buried in the view, instead of having to spell it out in every query. For example,

```
CREATE VIEW myremote_pg_proc AS
  SELECT *
    FROM dblink('dbname=postgres options=-csearch_path=',
                'select proname, prosrc from pg_proc')
    AS t1(proname name, prosrc text);

SELECT * FROM myremote_pg_proc WHERE proname LIKE 'bytea%';
```

Examples

```
SELECT * FROM dblink('dbname=postgres options=-csearch_path=',
                    'select proname, prosrc from pg_proc')
  AS t1(proname name, prosrc text) WHERE proname LIKE 'bytea%';
 proname | prosrc
-----+-----
 byteacat | byteacat
 byteaeq  | byteaeq
 bytealt  | bytealt
 byteale  | byteale
 byteagt  | byteagt
 byteage  | byteage
 byteane  | byteane
 byteacmp | byteacmp
 bytealike | bytealike
 byteanlike | byteanlike
 byteain  | byteain
 byteaout | byteaout
(12 rows)
```

```
SELECT dblink_connect('dbname=postgres options=-csearch_path=');
 dblink_connect
-----
      OK
(1 row)
```

```
SELECT * FROM dblink('select proname, prosrc from pg_proc')
  AS t1(proname name, prosrc text) WHERE proname LIKE 'bytea%';
 proname | prosrc
-----+-----
 byteacat | byteacat
 byteaeq  | byteaeq
 bytealt  | bytealt
 byteale  | byteale
 byteagt  | byteagt
 byteage  | byteage
 byteane  | byteane
 byteacmp | byteacmp
 bytealike | bytealike
 byteanlike | byteanlike
 byteain  | byteain
 byteaout | byteaout
(12 rows)
```

```
SELECT dblink_connect('myconn', 'dbname=regression options=-csearch_path=');
```


Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

dblink_connect

OK

(1 row)

```
SELECT * FROM dblink('myconn', 'select proname, prosrc from pg_proc')
  AS t1(proname name, prosrc text) WHERE proname LIKE 'bytea%';
```

proname		prosrc
-----+-----		
bytearecv		bytearecv
byteasend		byteasend
byteale		byteale
byteagt		byteagt
byteage		byteage
byteane		byteane
byteacmp		byteacmp
bytealike		bytealike
byteanlike		byteanlike
byteacat		byteacat
byteaeq		byteaeq
bytealt		bytealt
byteain		byteain
byteaout		byteaout

(14 rows)

dblink_exec

`dblink_exec` — executes a command in a remote database

Synopsis

```
dblink_exec(text connname, text sql [, bool fail_on_error]) returns text
dblink_exec(text connstr, text sql [, bool fail_on_error]) returns text
dblink_exec(text sql [, bool fail_on_error]) returns text
```

Description

`dblink_exec` executes a command (that is, any SQL statement that doesn't return rows) in a remote database.

When two `text` arguments are given, the first one is first looked up as a persistent connection's name; if found, the command is executed on that connection. If not found, the first argument is treated as a connection info string as for `dblink_connect`, and the indicated connection is made just for the duration of this command.

Arguments

connname

Name of the connection to use; omit this parameter to use the unnamed connection.

connstr

A connection info string, as previously described for `dblink_connect`.

sql

The SQL command that you wish to execute in the remote database, for example `insert into foo values (0, 'a', '{"a0","b0","c0"}')`.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a NOTICE, and the function's return value is set to `ERROR`.

Return Value

Returns status, either the command's status string or `ERROR`.

Examples

```
SELECT dblink_connect('dbname=dblink_test_standby');
 dblink_connect
-----
 OK
(1 row)

SELECT dblink_exec('insert into foo values(21, ''z'', ''{"a0","b0","c0"}'');');
 dblink_exec
-----
 INSERT 943366 1
(1 row)

SELECT dblink_connect('myconn', 'dbname=regression');
 dblink_connect
```

```
-----
OK
(1 row)

SELECT dblink_exec('myconn', 'insert into foo values(21, ''z'',
''{"a0","b0","c0"}'');');
      dblink_exec
-----

INSERT 6432584 1
(1 row)

SELECT dblink_exec('myconn', 'insert into pg_class values (''foo''),false);
NOTICE:  sql error
DETAIL:  ERROR:  null value in column "relnamespace" violates not-null constraint

      dblink_exec
-----

ERROR
(1 row)
```

dblink_open

`dblink_open` — opens a cursor in a remote database

Synopsis

```
dblink_open(text cursorname, text sql [, bool fail_on_error]) returns text
dblink_open(text connname, text cursorname, text sql [, bool fail_on_error]) returns
text
```

Description

`dblink_open()` opens a cursor in a remote database. The cursor can subsequently be manipulated with `dblink_fetch()` and `dblink_close()`.

Arguments

connname

Name of the connection to use; omit this parameter to use the unnamed connection.

cursorname

The name to assign to this cursor.

sql

The `SELECT` statement that you wish to execute in the remote database, for example `select * from pg_class`.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a NOTICE, and the function's return value is set to `ERROR`.

Return Value

Returns status, either `OK` or `ERROR`.

Notes

Since a cursor can only persist within a transaction, `dblink_open` starts an explicit transaction block (`BEGIN`) on the remote side, if the remote side was not already within a transaction. This transaction will be closed again when the matching `dblink_close` is executed. Note that if you use `dblink_exec` to change data between `dblink_open` and `dblink_close`, and then an error occurs or you use `dblink_disconnect` before `dblink_close`, your change *will be lost* because the transaction will be aborted.

Examples

```
SELECT dblink_connect('dbname=postgres options=-csearch_path=');
dblink_connect
-----
OK
(1 row)

SELECT dblink_open('foo', 'select proname, prosrc from pg_proc');
dblink_open
-----
OK
(1 row)
```

dblink_fetch

`dblink_fetch` — returns rows from an open cursor in a remote database

Synopsis

```
dblink_fetch(text cursorname, int howmany [, bool fail_on_error]) returns setof record
dblink_fetch(text connname, text cursorname, int howmany [, bool fail_on_error])
returns setof record
```

Description

`dblink_fetch` fetches rows from a cursor previously established by `dblink_open`.

Arguments

connname

Name of the connection to use; omit this parameter to use the unnamed connection.

cursorname

The name of the cursor to fetch from.

howmany

The maximum number of rows to retrieve. The next *howmany* rows are fetched, starting at the current cursor position, moving forward. Once the cursor has reached its end, no more rows are produced.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a NOTICE, and the function returns no rows.

Return Value

The function returns the row(s) fetched from the cursor. To use this function, you will need to specify the expected set of columns, as previously discussed for `dblink`.

Notes

On a mismatch between the number of return columns specified in the `FROM` clause, and the actual number of columns returned by the remote cursor, an error will be thrown. In this event, the remote cursor is still advanced by as many rows as it would have been if the error had not occurred. The same is true for any other error occurring in the local query after the remote `FETCH` has been done.

Examples

```
SELECT dblink_connect('dbname=postgres options=-csearch_path=');
dblink_connect
```

```
-----
OK
(1 row)
```

```
SELECT dblink_open('foo', 'select proname, prosrc from pg_proc where proname like
''bytea%'');
dblink_open
```

```
-----
OK
(1 row)
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

```
SELECT * FROM dblink_fetch('foo', 5) AS (funcname name, source text);
 funcname | source
-----+-----
 byteacat | byteacat
 byteacmp | byteacmp
 byteaeq  | byteaeq
 byteage  | byteage
 byteagt  | byteagt
(5 rows)
```

```
SELECT * FROM dblink_fetch('foo', 5) AS (funcname name, source text);
 funcname | source
-----+-----
 byteain  | byteain
 byteale  | byteale
 bytealike| bytealike
 bytealt  | bytealt
 byteane  | byteane
(5 rows)
```

```
SELECT * FROM dblink_fetch('foo', 5) AS (funcname name, source text);
 funcname | source
-----+-----
 byteanlike| byteanlike
 byteaout  | byteaout
(2 rows)
```

```
SELECT * FROM dblink_fetch('foo', 5) AS (funcname name, source text);
 funcname | source
-----+-----
(0 rows)
```

dblink_close

`dblink_close` — closes a cursor in a remote database

Synopsis

```
dblink_close(text cursorname [, bool fail_on_error]) returns text
dblink_close(text connname, text cursorname [, bool fail_on_error]) returns text
```

Description

`dblink_close` closes a cursor previously opened with `dblink_open`.

Arguments

connname

Name of the connection to use; omit this parameter to use the unnamed connection.

cursorname

The name of the cursor to close.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a NOTICE, and the function's return value is set to `ERROR`.

Return Value

Returns status, either `OK` or `ERROR`.

Notes

If `dblink_open` started an explicit transaction block, and this is the last remaining open cursor in this connection, `dblink_close` will issue the matching `COMMIT`.

Examples

```
SELECT dblink_connect('dbname=postgres options=-csearch_path=');
 dblink_connect
-----
 OK
(1 row)

SELECT dblink_open('foo', 'select proname, prosrc from pg_proc');
 dblink_open
-----
 OK
(1 row)

SELECT dblink_close('foo');
 dblink_close
-----
 OK
(1 row)
```

dblink_get_connections

`dblink_get_connections` — returns the names of all open named dblink connections

Synopsis

```
dblink_get_connections() returns text[]
```

Description

`dblink_get_connections` returns an array of the names of all open named dblink connections.

Return Value

Returns a text array of connection names, or NULL if none.

Examples

```
SELECT dblink_get_connections();
```


dblink_error_message

`dblink_error_message` — gets last error message on the named connection

Synopsis

```
dblink_error_message(text connname) returns text
```

Description

`dblink_error_message` fetches the most recent remote error message for a given connection.

Arguments

connname

Name of the connection to use.

Return Value

Returns last error message, or OK if there has been no error in this connection.

Notes

When asynchronous queries are initiated by `dblink_send_query`, the error message associated with the connection might not get updated until the server's response message is consumed. This typically means that `dblink_is_busy` or `dblink_get_result` should be called prior to `dblink_error_message`, so that any error generated by the asynchronous query will be visible.

Examples

```
SELECT dblink_error_message('dtest1');
```

dblink_send_query

`dblink_send_query` — sends an async query to a remote database

Synopsis

```
dblink_send_query(text connname, text sql) returns int
```

Description

`dblink_send_query` sends a query to be executed asynchronously, that is, without immediately waiting for the result. There must not be an async query already in progress on the connection.

After successfully dispatching an async query, completion status can be checked with `dblink_is_busy`, and the results are ultimately collected with `dblink_get_result`. It is also possible to attempt to cancel an active async query using `dblink_cancel_query`.

Arguments

connname

Name of the connection to use.

sql

The SQL statement that you wish to execute in the remote database, for example `select * from pg_class`.

Return Value

Returns 1 if the query was successfully dispatched, 0 otherwise.

Examples

```
SELECT dblink_send_query('dtest1', 'SELECT * FROM foo WHERE f1 < 3');
```

dblink_is_busy

`dblink_is_busy` — checks if connection is busy with an async query

Synopsis

```
dblink_is_busy(text connname) returns int
```

Description

`dblink_is_busy` tests whether an async query is in progress.

Arguments

connname

Name of the connection to check.

Return Value

Returns 1 if connection is busy, 0 if it is not busy. If this function returns 0, it is guaranteed that `dblink_get_result` will not block.

Examples

```
SELECT dblink_is_busy('dtest1');
```

dblink_get_notify

dblink_get_notify — retrieve async notifications on a connection

Synopsis

```
dblink_get_notify() returns setof (notify_name text, be_pid int, extra text)
dblink_get_notify(text connname) returns setof (notify_name text, be_pid int, extra
text)
```

Description

dblink_get_notify retrieves notifications on either the unnamed connection, or on a named connection if specified. To receive notifications via dblink, LISTEN must first be issued, using dblink_exec. For details see [LISTEN](#) and [NOTIFY](#).

Arguments

connname

The name of a named connection to get notifications on.

Return Value

Returns setof (notify_name text, be_pid int, extra text), or an empty set if none.

Examples

```
SELECT dblink_exec('LISTEN virtual');
 dblink_exec
-----
LISTEN
(1 row)
```

```
SELECT * FROM dblink_get_notify();
 notify_name | be_pid | extra
-----+-----+-----
(0 rows)
```

```
NOTIFY virtual;
NOTIFY
```

```
SELECT * FROM dblink_get_notify();
 notify_name | be_pid | extra
-----+-----+-----
virtual      |    1229 |
(1 row)
```

dblink_get_result

`dblink_get_result` — gets an async query result

Synopsis

`dblink_get_result(text connname [, bool fail_on_error])` returns setof record

Description

`dblink_get_result` collects the results of an asynchronous query previously sent with `dblink_send_query`. If the query is not already completed, `dblink_get_result` will wait until it is.

Arguments

connname

Name of the connection to use.

fail_on_error

If true (the default when omitted) then an error thrown on the remote side of the connection causes an error to also be thrown locally. If false, the remote error is locally reported as a NOTICE, and the function returns no rows.

Return Value

For an async query (that is, an SQL statement returning rows), the function returns the row(s) produced by the query. To use this function, you will need to specify the expected set of columns, as previously discussed for `dblink`.

For an async command (that is, an SQL statement not returning rows), the function returns a single row with a single text column containing the command's status string. It is still necessary to specify that the result will have a single text column in the calling `FROM` clause.

Notes

This function *must* be called if `dblink_send_query` returned 1. It must be called once for each query sent, and one additional time to obtain an empty set result, before the connection can be used again.

When using `dblink_send_query` and `dblink_get_result`, `dblink` fetches the entire remote query result before returning any of it to the local query processor. If the query returns a large number of rows, this can result in transient memory bloat in the local session. It may be better to open such a query as a cursor with `dblink_open` and then fetch a manageable number of rows at a time. Alternatively, use plain `dblink()`, which avoids memory bloat by spooling large result sets to disk.

Examples

```
contrib_regression=# SELECT dblink_connect('dtest1', 'dbname=contrib_regression');
dblink_connect
-----
OK
(1 row)

contrib_regression=# SELECT * FROM
contrib_regression-# dblink_send_query('dtest1', 'select * from foo where f1 < 3') AS
t1;
t1
----
1
(1 row)
```

```
contrib_regression=# SELECT * FROM dblink_get_result('dtest1') AS t1(f1 int, f2 text,
f3 text[]);
 f1 | f2 |      f3
-----+-----+-----
  0 | a  | {a0,b0,c0}
  1 | b  | {a1,b1,c1}
  2 | c  | {a2,b2,c2}
(3 rows)
```

```
contrib_regression=# SELECT * FROM dblink_get_result('dtest1') AS t1(f1 int, f2 text,
f3 text[]);
 f1 | f2 | f3
-----+-----+-----
(0 rows)
```

```
contrib_regression=# SELECT * FROM
contrib_regression=# dblink_send_query('dtest1', 'select * from foo where f1 < 3;
select * from foo where f1 > 6') AS t1;
t1
----
  1
(1 row)
```

```
contrib_regression=# SELECT * FROM dblink_get_result('dtest1') AS t1(f1 int, f2 text,
f3 text[]);
 f1 | f2 |      f3
-----+-----+-----
  0 | a  | {a0,b0,c0}
  1 | b  | {a1,b1,c1}
  2 | c  | {a2,b2,c2}
(3 rows)
```

```
contrib_regression=# SELECT * FROM dblink_get_result('dtest1') AS t1(f1 int, f2 text,
f3 text[]);
 f1 | f2 |      f3
-----+-----+-----
  7 | h  | {a7,b7,c7}
  8 | i  | {a8,b8,c8}
  9 | j  | {a9,b9,c9}
 10 | k  | {a10,b10,c10}
(4 rows)
```

```
contrib_regression=# SELECT * FROM dblink_get_result('dtest1') AS t1(f1 int, f2 text,
f3 text[]);
 f1 | f2 | f3
-----+-----+-----
(0 rows)
```

dblink_cancel_query

`dblink_cancel_query` — cancels any active query on the named connection

Synopsis

`dblink_cancel_query(text connname)` returns text

Description

`dblink_cancel_query` attempts to cancel any query that is in progress on the named connection. Note that this is not certain to succeed (since, for example, the remote query might already have finished). A cancel request simply improves the odds that the query will fail soon. You must still complete the normal query protocol, for example by calling `dblink_get_result`.

Arguments

connname

Name of the connection to use.

Return Value

Returns `OK` if the cancel request has been sent, or the text of an error message on failure.

Examples

```
SELECT dblink_cancel_query('dtest1');
```

dblink_get_pkey

`dblink_get_pkey` — returns the positions and field names of a relation's primary key fields

Synopsis

`dblink_get_pkey(text relname)` returns `setof dblink_pkey_results`

Description

`dblink_get_pkey` provides information about the primary key of a relation in the local database. This is sometimes useful in generating queries to be sent to remote databases.

Arguments

relname

Name of a local relation, for example `foo` or `myschema.mytab`. Include double quotes if the name is mixed-case or contains special characters, for example `"FooBar"`; without quotes, the string will be folded to lower case.

Return Value

Returns one row for each primary key field, or no rows if the relation has no primary key. The result row type is defined as

```
CREATE TYPE dblink_pkey_results AS (position int, colname text);
```

The `position` column simply runs from 1 to *N*; it is the number of the field within the primary key, not the number within the table's columns.

Examples

```
CREATE TABLE foobar (  
    f1 int,  
    f2 int,  
    f3 int,  
    PRIMARY KEY (f1, f2, f3)  
);  
CREATE TABLE  
  
SELECT * FROM dblink_get_pkey('foobar');  
 position | colname  
-----+-----  
        1 | f1  
        2 | f2  
        3 | f3  
(3 rows)
```


dblink_build_sql_insert

`dblink_build_sql_insert` — builds an `INSERT` statement using a local tuple, replacing the primary key field values with alternative supplied values

Synopsis

```
dblink_build_sql_insert(text relname,  
                        int2vector primary_key_attnums,  
                        integer num_primary_key_atts,  
                        text[] src_pk_att_vals_array,  
                        text[] tgt_pk_att_vals_array) returns text
```

Description

`dblink_build_sql_insert` can be useful in doing selective replication of a local table to a remote database. It selects a row from the local table based on primary key, and then builds an SQL `INSERT` command that will duplicate that row, but with the primary key values replaced by the values in the last argument. (To make an exact copy of the row, just specify the same values for the last two arguments.)

Arguments

relname

Name of a local relation, for example `foo` or `myschema.mytab`. Include double quotes if the name is mixed-case or contains special characters, for example `"FooBar"`; without quotes, the string will be folded to lower case.

primary_key_attnums

Attribute numbers (1-based) of the primary key fields, for example `1 2`.

num_primary_key_atts

The number of primary key fields.

src_pk_att_vals_array

Values of the primary key fields to be used to look up the local tuple. Each field is represented in text form. An error is thrown if there is no local row with these primary key values.

tgt_pk_att_vals_array

Values of the primary key fields to be placed in the resulting `INSERT` command. Each field is represented in text form.

Return Value

Returns the requested SQL statement as text.

Notes

As of PostgreSQL 9.0, the attribute numbers in *primary_key_attnums* are interpreted as logical column numbers, corresponding to the column's position in `SELECT * FROM relname`. Previous versions interpreted the numbers as physical column positions. There is a difference if any column(s) to the left of the indicated column have been dropped during the lifetime of the table.

Examples

```
SELECT dblink_build_sql_insert('foo', '1 2', 2, '{"1", "a"}', '{"1", "b''a"}');  
       dblink_build_sql_insert
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

```
INSERT INTO foo(f1,f2,f3) VALUES('1','b' 'a','1')
(1 row)
```

dblink_build_sql_delete

`dblink_build_sql_delete` — builds a DELETE statement using supplied values for primary key field values

Synopsis

```
dblink_build_sql_delete(text relname,  
                        int2vector primary_key_attnums,  
                        integer num_primary_key_atts,  
                        text[] tgt_pk_att_vals_array) returns text
```

Description

`dblink_build_sql_delete` can be useful in doing selective replication of a local table to a remote database. It builds an SQL DELETE command that will delete the row with the given primary key values.

Arguments

relname

Name of a local relation, for example `foo` or `myschema.mytab`. Include double quotes if the name is mixed-case or contains special characters, for example `"FooBar"`; without quotes, the string will be folded to lower case.

primary_key_attnums

Attribute numbers (1-based) of the primary key fields, for example `1 2`.

num_primary_key_atts

The number of primary key fields.

tgt_pk_att_vals_array

Values of the primary key fields to be used in the resulting DELETE command. Each field is represented in text form.

Return Value

Returns the requested SQL statement as text.

Notes

As of PostgreSQL 9.0, the attribute numbers in *primary_key_attnums* are interpreted as logical column numbers, corresponding to the column's position in `SELECT * FROM relname`. Previous versions interpreted the numbers as physical column positions. There is a difference if any column(s) to the left of the indicated column have been dropped during the lifetime of the table.

Examples

```
SELECT dblink_build_sql_delete('MyFoo', '1 2', 2, '{"1", "b"}');  
      dblink_build_sql_delete  
-----  
DELETE FROM "MyFoo" WHERE f1='1' AND f2='b'  
(1 row)
```

dblink_build_sql_update

`dblink_build_sql_update` — builds an `UPDATE` statement using a local tuple, replacing the primary key field values with alternative supplied values

Synopsis

```
dblink_build_sql_update(text relname,  
                        int2vector primary_key_attnums,  
                        integer num_primary_key_atts,  
                        text[] src_pk_att_vals_array,  
                        text[] tgt_pk_att_vals_array) returns text
```

Description

`dblink_build_sql_update` can be useful in doing selective replication of a local table to a remote database. It selects a row from the local table based on primary key, and then builds an SQL `UPDATE` command that will duplicate that row, but with the primary key values replaced by the values in the last argument. (To make an exact copy of the row, just specify the same values for the last two arguments.) The `UPDATE` command always assigns all fields of the row — the main difference between this and `dblink_build_sql_insert` is that it's assumed that the target row already exists in the remote table.

Arguments

relname

Name of a local relation, for example `foo` or `myschema.mytab`. Include double quotes if the name is mixed-case or contains special characters, for example `"FooBar"`; without quotes, the string will be folded to lower case.

primary_key_attnums

Attribute numbers (1-based) of the primary key fields, for example `1 2`.

num_primary_key_atts

The number of primary key fields.

src_pk_att_vals_array

Values of the primary key fields to be used to look up the local tuple. Each field is represented in text form. An error is thrown if there is no local row with these primary key values.

tgt_pk_att_vals_array

Values of the primary key fields to be placed in the resulting `UPDATE` command. Each field is represented in text form.

Return Value

Returns the requested SQL statement as text.

Notes

As of PostgreSQL 9.0, the attribute numbers in *primary_key_attnums* are interpreted as logical column numbers, corresponding to the column's position in `SELECT * FROM relname`. Previous versions interpreted the numbers as physical column positions. There is a difference if any column(s) to the left of the indicated column have been dropped during the lifetime of the table.

Examples

```
SELECT dblink_build_sql_update('foo', '1 2', 2, '{"1", "a"}', '{"1", "b"}');
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

dblink_build_sql_update

```
UPDATE foo SET f1='1',f2='b',f3='1' WHERE f1='1' AND f2='b'
(1 row)
```

F.5. pgcrypto — cryptographic functions

The `pgcrypto` module provides cryptographic functions for Postgres Pro.

This module is considered “trusted”, that is, it can be installed by non-superusers who have `CREATE` privilege on the current database.

`pgcrypto` requires OpenSSL and won't be installed if OpenSSL support was not selected when Postgres Pro was built.

F.5.1. General Hashing Functions

F.5.1.1. `digest()`

`digest(data text, type text)` returns `bytea`
`digest(data bytea, type text)` returns `bytea`

Computes a binary hash of the given *data*. *type* is the algorithm to use. Standard algorithms are `md5`, `sha1`, `sha224`, `sha256`, `sha384` and `sha512`. Moreover, any digest algorithm OpenSSL supports is automatically picked up.

If you want the digest as a hexadecimal string, use `encode()` on the result. For example:

```
CREATE OR REPLACE FUNCTION sha1(bytea) returns text AS $$  
    SELECT encode(digest($1, 'sha1'), 'hex')  
$$ LANGUAGE SQL STRICT IMMUTABLE;
```

F.5.1.2. `hmac()`

`hmac(data text, key text, type text)` returns `bytea`
`hmac(data bytea, key bytea, type text)` returns `bytea`

Calculates hashed MAC for *data* with key *key*. *type* is the same as in `digest()`.

This is similar to `digest()` but the hash can only be recalculated knowing the key. This prevents the scenario of someone altering data and also changing the hash to match.

If the key is larger than the hash block size it will first be hashed and the result will be used as key.

F.5.2. Password Hashing Functions

The functions `crypt()` and `gen_salt()` are specifically designed for hashing passwords. `crypt()` does the hashing and `gen_salt()` prepares algorithm parameters for it.

The algorithms in `crypt()` differ from the usual MD5 or SHA1 hashing algorithms in the following respects:

1. They are slow. As the amount of data is so small, this is the only way to make brute-forcing passwords hard.
2. They use a random value, called the *salt*, so that users having the same password will have different encrypted passwords. This is also an additional defense against reversing the algorithm.
3. They include the algorithm type in the result, so passwords hashed with different algorithms can co-exist.
4. Some of them are adaptive — that means when computers get faster, you can tune the algorithm to be slower, without introducing incompatibility with existing passwords.

[Table F.5](#) lists the algorithms supported by the `crypt()` function.

Table F.5. Supported Algorithms for `crypt()`

Algorithm	Max Password Length	Adaptive?	Salt Bits	Output Length	Description
bf	72	yes	128	60	Blowfish-based, variant 2a
md5	unlimited	no	48	34	MD5-based crypt
xdes	8	yes	24	20	Extended DES
des	8	no	12	13	Original UNIX crypt

F.5.2.1. `crypt()`

`crypt(password text, salt text)` returns text

Calculates a crypt(3)-style hash of *password*. When storing a new password, you need to use `gen_salt()` to generate a new *salt* value. To check a password, pass the stored hash value as *salt*, and test whether the result matches the stored value.

Example of setting a new password:

```
UPDATE ... SET pswhash = crypt('new password', gen_salt('md5'));
```

Example of authentication:

```
SELECT (pswhash = crypt('entered password', pswhash)) AS pswmatch FROM ... ;
```

This returns `true` if the entered password is correct.

F.5.2.2. `gen_salt()`

`gen_salt(type text [, iter_count integer])` returns text

Generates a new random salt string for use in `crypt()`. The salt string also tells `crypt()` which algorithm to use.

The *type* parameter specifies the hashing algorithm. The accepted types are: `des`, `xdes`, `md5` and `bf`.

The *iter_count* parameter lets the user specify the iteration count, for algorithms that have one. The higher the count, the more time it takes to hash the password and therefore the more time to break it. Although with too high a count the time to calculate a hash may be several years — which is somewhat impractical. If the *iter_count* parameter is omitted, the default iteration count is used. Allowed values for *iter_count* depend on the algorithm and are shown in [Table F.6](#).

Table F.6. Iteration Counts for `crypt()`

Algorithm	Default	Min	Max
xdes	725	1	16777215
bf	6	4	31

For `xdes` there is an additional limitation that the iteration count must be an odd number.

To pick an appropriate iteration count, consider that the original DES crypt was designed to have the speed of 4 hashes per second on the hardware of that time. Slower than 4 hashes per second would probably dampen usability. Faster than 100 hashes per second is probably too fast.

[Table F.7](#) gives an overview of the relative slowness of different hashing algorithms. The table shows how much time it would take to try all combinations of characters in an 8-character password, assuming that the password contains either only lower case letters, or upper- and lower-case letters and numbers. In the `crypt-bf` entries, the number after a slash is the *iter_count* parameter of `gen_salt`.

Table F.7. Hash Algorithm Speeds

Algorithm	Hashes/sec	For [a-z]	For [A-Za-z0-9]	Duration relative to md5 hash
crypt-bf/8	1792	4 years	3927 years	100k
crypt-bf/7	3648	2 years	1929 years	50k
crypt-bf/6	7168	1 year	982 years	25k
crypt-bf/5	13504	188 days	521 years	12.5k
crypt-md5	171584	15 days	41 years	1k
crypt-des	23221568	157.5 minutes	108 days	7
sha1	37774272	90 minutes	68 days	4
md5 (hash)	150085504	22.5 minutes	17 days	1

Notes:

- The machine used is an Intel Mobile Core i3.
- `crypt-des` and `crypt-md5` algorithm numbers are taken from John the Ripper v1.6.38 `-test` output.
- `md5 hash` numbers are from `mdcrack` 1.2.
- `sha1` numbers are from `lcrack-20031130-beta`.
- `crypt-bf` numbers are taken using a simple program that loops over 1000 8-character passwords. That way the speed with different numbers of iterations can be shown. For reference: `john -test` shows 13506 loops/sec for `crypt-bf/5`. (The very small difference in results is in accordance with the fact that the `crypt-bf` implementation in `pgcrypto` is the same one used in John the Ripper.)

Note that “try all combinations” is not a realistic exercise. Usually password cracking is done with the help of dictionaries, which contain both regular words and various mutations of them. So, even somewhat word-like passwords could be cracked much faster than the above numbers suggest, while a 6-character non-word-like password may escape cracking. Or not.

F.5.3. PGP Encryption Functions

The functions here implement the encryption part of the OpenPGP ([RFC 4880](#)) standard. Supported are both symmetric-key and public-key encryption.

An encrypted PGP message consists of 2 parts, or *packets*:

- Packet containing a session key — either symmetric-key or public-key encrypted.
- Packet containing data encrypted with the session key.

When encrypting with a symmetric key (i.e., a password):

1. The given password is hashed using a String2Key (S2K) algorithm. This is rather similar to `crypt()` algorithms — purposefully slow and with random salt — but it produces a full-length binary key.
2. If a separate session key is requested, a new random key will be generated. Otherwise the S2K key will be used directly as the session key.
3. If the S2K key is to be used directly, then only S2K settings will be put into the session key packet. Otherwise the session key will be encrypted with the S2K key and put into the session key packet.

When encrypting with a public key:

1. A new random session key is generated.
2. It is encrypted using the public key and put into the session key packet.

In either case the data to be encrypted is processed as follows:

1. Optional data-manipulation: compression, conversion to UTF-8, and/or conversion of line-endings.
2. The data is prefixed with a block of random bytes. This is equivalent to using a random IV.
3. A SHA1 hash of the random prefix and data is appended.
4. All this is encrypted with the session key and placed in the data packet.

F.5.3.1. `pgp_sym_encrypt()`

`pgp_sym_encrypt(data text, psw text [, options text])` returns `bytea`
`pgp_sym_encrypt_bytea(data bytea, psw text [, options text])` returns `bytea`

Encrypt *data* with a symmetric PGP key *psw*. The *options* parameter can contain option settings, as described below.

F.5.3.2. `pgp_sym_decrypt()`

`pgp_sym_decrypt(msg bytea, psw text [, options text])` returns `text`
`pgp_sym_decrypt_bytea(msg bytea, psw text [, options text])` returns `bytea`

Decrypt a symmetric-key-encrypted PGP message.

Decrypting `bytea` data with `pgp_sym_decrypt` is disallowed. This is to avoid outputting invalid character data. Decrypting originally textual data with `pgp_sym_decrypt_bytea` is fine.

The *options* parameter can contain option settings, as described below.

F.5.3.3. `pgp_pub_encrypt()`

`pgp_pub_encrypt(data text, key bytea [, options text])` returns `bytea`
`pgp_pub_encrypt_bytea(data bytea, key bytea [, options text])` returns `bytea`

Encrypt *data* with a public PGP key *key*. Giving this function a secret key will produce an error.

The *options* parameter can contain option settings, as described below.

F.5.3.4. `pgp_pub_decrypt()`

`pgp_pub_decrypt(msg bytea, key bytea [, psw text [, options text]])` returns `text`
`pgp_pub_decrypt_bytea(msg bytea, key bytea [, psw text [, options text]])` returns `bytea`

Decrypt a public-key-encrypted message. *key* must be the secret key corresponding to the public key that was used to encrypt. If the secret key is password-protected, you must give the password in *psw*. If there is no password, but you want to specify options, you need to give an empty password.

Decrypting `bytea` data with `pgp_pub_decrypt` is disallowed. This is to avoid outputting invalid character data. Decrypting originally textual data with `pgp_pub_decrypt_bytea` is fine.

The *options* parameter can contain option settings, as described below.

F.5.3.5. `pgp_key_id()`

`pgp_key_id(bytea)` returns `text`

`pgp_key_id` extracts the key ID of a PGP public or secret key. Or it gives the key ID that was used for encrypting the data, if given an encrypted message.

It can return 2 special key IDs:

- `SYMKEY`

The message is encrypted with a symmetric key.

- `ANYKEY`

The message is public-key encrypted, but the key ID has been removed. That means you will need to try all your secret keys on it to see which one decrypts it. `pgcrypto` itself does not produce such messages.

Note that different keys may have the same ID. This is rare but a normal event. The client application should then try to decrypt with each one, to see which fits — like handling `ANYKEY`.

F.5.3.6. `armor()`, `dearmor()`

```
armor(data bytea [ , keys text[], values text[] ]) returns text
dearmor(data text) returns bytea
```

These functions wrap/unwrap binary data into PGP ASCII-armor format, which is basically Base64 with CRC and additional formatting.

If the *keys* and *values* arrays are specified, an *armor header* is added to the armored format for each key/value pair. Both arrays must be single-dimensional, and they must be of the same length. The keys and values cannot contain any non-ASCII characters.

F.5.3.7. `pgp_armor_headers`

```
pgp_armor_headers(data text, key out text, value out text) returns setof record
```

`pgp_armor_headers()` extracts the armor headers from *data*. The return value is a set of rows with two columns, key and value. If the keys or values contain any non-ASCII characters, they are treated as UTF-8.

F.5.3.8. Options for PGP Functions

Options are named to be similar to GnuPG. An option's value should be given after an equal sign; separate options from each other with commas. For example:

```
pgp_sym_encrypt(data, psw, 'compress-algo=1, cipher-algo=aes256')
```

All of the options except `convert-crlf` apply only to encrypt functions. Decrypt functions get the parameters from the PGP data.

The most interesting options are probably `compress-algo` and `unicode-mode`. The rest should have reasonable defaults.

F.5.3.8.1. `cipher-algo`

Which cipher algorithm to use.

Values: `bf`, `aes128`, `aes192`, `aes256`, `3des`, `cast5`

Default: `aes128`

Applies to: `pgp_sym_encrypt`, `pgp_pub_encrypt`

F.5.3.8.2. `compress-algo`

Which compression algorithm to use. Only available if Postgres Pro was built with `zlib`.

Values:

0 - no compression

1 - ZIP compression

2 - ZLIB compression (= ZIP plus meta-data and block CRCs)

Default: 0

Applies to: `pgp_sym_encrypt`, `pgp_pub_encrypt`

F.5.3.8.3. `compress-level`

How much to compress. Higher levels compress smaller but are slower. 0 disables compression.

Values: 0, 1-9

Default: 6

Applies to: pgp_sym_encrypt, pgp_pub_encrypt

F.5.3.8.4. convert-crlf

Whether to convert `\n` into `\r\n` when encrypting and `\r\n` to `\n` when decrypting. RFC 4880 specifies that text data should be stored using `\r\n` line-feeds. Use this to get fully RFC-compliant behavior.

Values: 0, 1

Default: 0

Applies to: pgp_sym_encrypt, pgp_pub_encrypt, pgp_sym_decrypt, pgp_pub_decrypt

F.5.3.8.5. disable-mdc

Do not protect data with SHA-1. The only good reason to use this option is to achieve compatibility with ancient PGP products, predating the addition of SHA-1 protected packets to RFC 4880. Recent gnupg.org and pgp.com software supports it fine.

Values: 0, 1

Default: 0

Applies to: pgp_sym_encrypt, pgp_pub_encrypt

F.5.3.8.6. sess-key

Use separate session key. Public-key encryption always uses a separate session key; this option is for symmetric-key encryption, which by default uses the S2K key directly.

Values: 0, 1

Default: 0

Applies to: pgp_sym_encrypt

F.5.3.8.7. s2k-mode

Which S2K algorithm to use.

Values:

0 - Without salt. Dangerous!

1 - With salt but with fixed iteration count.

3 - Variable iteration count.

Default: 3

Applies to: pgp_sym_encrypt

F.5.3.8.8. s2k-count

The number of iterations of the S2K algorithm to use. It must be a value between 1024 and 65011712, inclusive.

Default: A random value between 65536 and 253952

Applies to: pgp_sym_encrypt, only with s2k-mode=3

F.5.3.8.9. s2k-digest-algo

Which digest algorithm to use in S2K calculation.

Values: md5, sha1

Default: sha1

Applies to: pgp_sym_encrypt

F.5.3.8.10. s2k-cipher-algo

Which cipher to use for encrypting separate session key.

Values: bf, aes, aes128, aes192, aes256

Default: use cipher-algo

Applies to: pgp_sym_encrypt

F.5.3.8.11. unicode-mode

Whether to convert textual data from database internal encoding to UTF-8 and back. If your database already is UTF-8, no conversion will be done, but the message will be tagged as UTF-8. Without this option it will not be.

Values: 0, 1

Default: 0

Applies to: pgp_sym_encrypt, pgp_pub_encrypt

F.5.3.9. Generating PGP Keys with GnuPG

To generate a new key:

```
gpg --gen-key
```

The preferred key type is “DSA and Elgamal”.

For RSA encryption you must create either DSA or RSA sign-only key as master and then add an RSA encryption subkey with `gpg --edit-key`.

To list keys:

```
gpg --list-secret-keys
```

To export a public key in ASCII-armor format:

```
gpg -a --export KEYID > public.key
```

To export a secret key in ASCII-armor format:

```
gpg -a --export-secret-keys KEYID > secret.key
```

You need to use `dearmor()` on these keys before giving them to the PGP functions. Or if you can handle binary data, you can drop `-a` from the command.

For more details see `man gpg`, *The GNU Privacy Handbook* and other documentation on <https://www.gnupg.org/>.

F.5.3.10. Limitations of PGP Code

- No support for signing. That also means that it is not checked whether the encryption subkey belongs to the master key.
- No support for encryption key as master key. As such practice is generally discouraged, this should not be a problem.
- No support for several subkeys. This may seem like a problem, as this is common practice. On the other hand, you should not use your regular GPG/PGP keys with `pgcrypto`, but create new ones, as the usage scenario is rather different.

F.5.4. Raw Encryption Functions

These functions only run a cipher over data; they don't have any advanced features of PGP encryption. Therefore they have some major problems:

1. They use user key directly as cipher key.
2. They don't provide any integrity checking, to see if the encrypted data was modified.
3. They expect that users manage all encryption parameters themselves, even IV.
4. They don't handle text.

So, with the introduction of PGP encryption, usage of raw encryption functions is discouraged.

`encrypt(data bytea, key bytea, type text)` returns `bytea`
`decrypt(data bytea, key bytea, type text)` returns `bytea`

`encrypt_iv(data bytea, key bytea, iv bytea, type text)` returns `bytea`
`decrypt_iv(data bytea, key bytea, iv bytea, type text)` returns `bytea`

Encrypt/decrypt data using the cipher method specified by *type*. The syntax of the *type* string is:

`algorithm [- mode] [/pad: padding]`

where *algorithm* is one of:

- `bf` — Blowfish
- `aes` — AES (Rijndael-128, -192 or -256)

and *mode* is one of:

- `cbc` — next block depends on previous (default)
- `ecb` — each block is encrypted separately (for testing only)

and *padding* is one of:

- `pkcs` — data may be any length (default)
- `none` — data must be multiple of cipher block size

So, for example, these are equivalent:

```
encrypt(data, 'fooz', 'bf')
encrypt(data, 'fooz', 'bf-cbc/pad:pkcs')
```

In `encrypt_iv` and `decrypt_iv`, the *iv* parameter is the initial value for the CBC mode; it is ignored for ECB. It is clipped or padded with zeroes if not exactly block size. It defaults to all zeroes in the functions without this parameter.

F.5.5. Random-Data Functions

`gen_random_bytes(count integer)` returns `bytea`

Returns *count* cryptographically strong random bytes. At most 1024 bytes can be extracted at a time. This is to avoid draining the randomness generator pool.

`gen_random_uuid()` returns `uuid`

Returns a version 4 (random) UUID. (Obsolete, this function internally calls the [core function](#) of the same name.)

F.5.6. Notes

F.5.6.1. Configuration

`pgcrypto` configures itself according to the findings of the main Postgres Pro `configure` script. The options that affect it are `--with-zlib` and `--with-ssl=openssl`.

When compiled with `zlib`, PGP encryption functions are able to compress data before encrypting.

`pgcrypto` requires `OpenSSL`. Otherwise, it will not be built or installed.

When compiled against `OpenSSL 3.0.0` and later versions, the legacy provider must be activated in the `openssl.cnf` configuration file in order to use older ciphers like DES or Blowfish.

F.5.6.2. NULL Handling

As is standard in SQL, all functions return `NULL`, if any of the arguments are `NULL`. This may create security risks on careless usage.

F.5.6.3. Security Limitations

All `pgcrypto` functions run inside the database server. That means that all the data and passwords move between `pgcrypto` and client applications in clear text. Thus you must:

1. Connect locally or use SSL connections.
2. Trust both system and database administrator.

If you cannot, then better do crypto inside client application.

The implementation does not resist *side-channel attacks*. For example, the time required for a `pgcrypto` decryption function to complete varies among ciphertexts of a given size.

F.5.7. Author

Marko Kreen <markokr@gmail.com>

`pgcrypto` uses code from the following sources:

Algorithm	Author	Source origin
DES crypt	David Burren and others	FreeBSD libcrypt
MD5 crypt	Poul-Henning Kamp	FreeBSD libcrypt
Blowfish crypt	Solar Designer	www.openwall.com

F.6. pg_proaudit — enables detailed logging of various security events

The `pg_proaudit` extension enables detailed logging of various security events.

`pg_proaudit` works in parallel with the standard PostgreSQL logging solutions (logging collector) and does not depend on them. Security event log of the `pg_proaudit` extension is stored separately from the server log. At the Postgres Pro Shardman startup, `pg_proaudit` launches a special background process to log security events.

Logging rules are stored in the `pg_proaudit.conf` configuration file located in the cluster data directory (PGDATA). It is a text file that can be edited directly using operating system facilities. To modify the file using SQL, you can use several `pg_proaudit` functions. The `pg_proaudit_settings` view displays the current `pg_proaudit` rules, even if they have not been saved into the `pg_proaudit.conf` file yet.

All the logged events belong to the following classes:

- DDL commands for creating, changing, and deleting DBMS objects (databases, tablespaces, schemas, tables, views, sequences, languages, functions)
- access control commands for database objects (GRANT, REVOKE)
- DML commands for access to database objects (INSERT, UPDATE, DELETE, SELECT, TRUNCATE for tables and/or views, EXECUTE for functions)
- database authentication/disconnection events
- all commands executed by a particular user

Security events can be logged both in the centralized logging solution of the operating system (syslog) and in the standard file-system files. Event logs can be written both into the syslog and into the files simultaneously. For clear identification, all `pg_proaudit` records in the syslog are marked with `AUDIT`. Event log files are written in the CSV format. Each event is logged on a separate line that contains the following fields:

- date and time of the event
- current user name
- database name
- server process ID (PID)
- severity level: `INFO` or `ERROR`
- serial number of the command in a session
- subcommand number in complex commands (CREATE TABLE ... AS SELECT ...)
- operator name
- object type
- object name
- operator execution results: `SUCCESS` or `FAILURE`
- additional information like error message in case of `FAILURE` or connection parameters for the `AUTHENTICATE` event; unfinished authentication attempts are marked as [EOF. No credentials provided]
- text of the SQL command
- parameters of the command (for example, for `PREPARE`)
- session user name
- unique ID of the event (in the [UUIDv7](#) format)
- transaction ID (XID)
- virtual transaction ID (VXID)

An event contains the information about both `session_user` and `current_user` of the session. Therefore, it will be possible to identify the user, even if they execute the `SET ROLE` command to change user identifier.

The XID and VXID values can be zero if an event does not relate to a transaction (e. g. the `DISCONNECT` event).

You can define a directory to store security log files and set up log file rotation. `pg_proaudit` can switch to a new log file either after the specified time interval, or when the specified size of the log file is exceeded. This enables you to define a workflow for cleaning up security event logs.

Postgres Pro user with the `SUPERUSER` attribute should grant access to the `pg_proaudit` extension and security event log files only to the user with the information security administrator role.

F.6.1. Installing the `pg_proaudit` Extension

The `pg_proaudit` extension is a built-in extension included into Postgres Pro Shardman. To enable `pg_proaudit`, complete the following steps:

1. Add `pg_proaudit` to the `shared_preload_libraries` variable in the `postgresql.conf` file:

```
shared_preload_libraries = 'pg_proaudit'
```

2. Restart the database server for the changes to take effect. To verify that the `pg_proaudit` library was installed correctly, you can run the following command:

```
SHOW shared_preload_libraries;
```

3. In each database for which you are going to log security events, create the `pg_proaudit` extension using the following query:

```
CREATE EXTENSION pg_proaudit;
```

The `pg_proaudit` extension adds several functions for managing the `pg_proaudit.conf` file, the `pg_proaudit_settings` view that displays the current `pg_proaudit` rules and event triggers.

F.6.2. Uninstalling the `pg_proaudit` Extension

To properly uninstall `pg_proaudit`, complete the following steps:

1. Delete the `pg_proaudit` extension using the following query:

```
DROP EXTENSION pg_proaudit;
```

2. Remove `pg_proaudit` from the `shared_preload_libraries` variable in the `postgresql.conf` file.

Skip this step if you have several databases in the cluster and you want to remove the extension only for one of them. In this case, it is recommended to remove logging rules related to the corresponding database prior to uninstalling the extension.

F.6.3. Functions to Configure Security Event Logging

To configure security event logging, `pg_proaudit` provides an SQL interface that consists of several functions and the `pg_proaudit_settings` view.

```
pg_proaudit_set_rule(db_name text, event_type text, object_type text, object_name text,  
                    role_name text, comment text)
```

Creates the logging rule with the specified parameters. When the `pg_proaudit_set_rule()` function completes, security event logging starts immediately, but the `pg_proaudit.conf` file is not updated. To save the changes in the `pg_proaudit.conf` file, call the `pg_proaudit_save()` function.

Arguments:

- *db_name* — name of the database for which the logging rule is established. An empty string or NULL specified in this argument means that events are logged for all databases where the `pg_proaudit` extension is created. When set to [current_database\(\)](#), events for the current database are logged.
- *event_type* — type of the event that needs to be logged, including SQL operator names, as well as `AUTHENTICATE` and `DISCONNECT` events. When set to `ALL`, as well as when an empty string or NULL is specified, enables logging for all events available for the specified object type. For example, for

the `TABLE` object type, the `ALL` keyword enables logging for commands `SELECT`, `INSERT`, `UPDATE`, `DELETE`, `TRUNCATE`, `COPY`, as well as `CREATE`, `ALTER`, `DROP`. You can also set up logging of event classes by specifying the following values: `ALL_DDL`, `ALL_DDL_NONTEMP`, `ALL_DML`, `ALL_DML_NONTEMP`, `ALL_MOD`, `ALL_PROC`, and `ALL_ROLE`. For the full list of possible `event_type` values, see [Section F.6.5](#).

- `object_type` — type of the object for which security events need to be logged. When set to `ALL`, as well as when an empty string or the `NULL` is specified, enables logging of events for all object types. For example, specify `FOREIGN TABLE` object type for the `SELECT` event to log all attempts to access foreign tables. Use `NULL` if `event_type` is set to `AUTHENTICATE`, `DISCONNECT`, `SET`, or `RESET`, and the `ROLE` value for all events that reference user actions, such as `CREATE USER` or `DROP USER`. The following object types are supported: `COMPOSITE TYPE`, `DATABASE`, `EVENT TRIGGER`, `FUNCTION`, `INDEX`, `PREPARED STATEMENT`, `ROLE`, `SEQUENCE`, `SCHEMA`, `TABLE`, `FOREIGN TABLE`, `TOAST TABLE`, `TABLESPACE`, `VIEW`, `MATERIALIZED VIEW`, `CATALOG RELATION`, and `CATALOG FUNCTION`.

Note

To log security events for system catalog objects, it is required to enable the `pg_proaudit.log_catalog_access` configuration parameter. Otherwise, these events will not be logged, even if the `CATALOG RELATION`, `CATALOG FUNCTION`, or `ALL` object type is specified in the rule.

- `object_name` — name of the object for which the logging rule is established. When an empty string or `NULL` is specified, enables logging of events for all object names.
- `role_name` — name of the role for which the logging rule is established. If specified, allows logging the actions caused by the user who has the privileges of the role. This means that at least one of the user session attributes, `session_user` or `current_user`, should either be equal to the `role_name`, or be directly or indirectly a member of this role. When an empty string or `NULL` is specified, enables logging of actions caused by any user. When set to `current_role`, the logging rule is established for the user current role.
- `comment` — comment to describe the created logging rule. This argument does not affect the rule execution and is not reflected in the log.

```
pg_proaudit_remove_rule(db_name text, event_type text, object_type text, object_name  
text, role_name text)
```

Removes the specific logging rule with the set parameters. To save the changes in the `pg_proaudit.conf` file, call the `pg_proaudit_save()` function.

Note

If the rule that you want to remove is configured to log the `DISCONNECT` event type, then disconnect events may still be logged after the rule is removed. For additional information, see [Section F.6.3.1](#).

Arguments:

- `db_name` — name of the database for which the logging rule needs to be removed.
- `event_type` — type of the event for which the logging rule needs to be removed. For the full list of possible `event_type` values, see [Section F.6.5](#).
- `object_type` — type of the object for which the logging rule needs to be removed.
- `object_name` — name of the object for which the logging rule needs to be removed.
- `role_name` — name of the role for which the logging rule needs to be removed.

```
pg_proaudit_show()
```

Returns logged events in a table view. This function is used by the `pg_proaudit_settings` view.

```
pg_proaudit_reload()
```

Reads logging configuration from the `pg_proaudit.conf` file. You must call this function if the `pg_proaudit.conf` file was modified manually using the operating system facilities.

`pg_proaudit_reset()`

Removes all logging rules. To save information about the canceled logging in the `pg_proaudit.conf` file, call the `pg_proaudit_save()` function.

`pg_proaudit_save()`

Saves logging rules from memory into the `pg_proaudit.conf` file. The `pg_proaudit.conf` file is located in the cluster data directory (PGDATA). You cannot change the `pg_proaudit.conf` file location.

F.6.3.1. Handling Database Disconnection Events

When a logging rule is configured for `DISCONNECT` event type, and a user authenticates in a database, then the `pg_proaudit` extension checks whether the corresponding disconnection event satisfies the rule and must be logged.

If the event must be logged, then the removal of the logging rule has no effect on the event. This event will be logged no matter if the logging rule exists at the moment of disconnect.

Consider the following example:

1. An administrator configures the following logging rule:

```
SELECT pg_proaudit_set_rule ('postgres', 'DISCONNECT', null, null, null, 'Any
disconnect from the postgres DB');
```

2. A user connects to the `postgres` database and goes through authentication.

For the rule in effect, a corresponding disconnection event will be logged later.

3. The administrator removes the logging rule, while the user is still connected to the database:

```
SELECT pg_proaudit_remove_rule('postgres', 'DISCONNECT', null, null, null);
```

4. The user disconnects from the database.

The corresponding disconnection event is logged despite the removal of the logging rule.

F.6.3.2. Handling Object Names

Comparison of string representations of names for database, user, role and other objects is case-insensitive even if the object name in the database is enclosed in quotes. These string representations are stored as lower-case strings.

Strings representing object names need to be passed to functions `pg_proaudit_set_rule` and `pg_proaudit_remove_rule` as regular strings: a single quote needs repeating, and other characters are passed as is, without escaping or wrapping.

For example, rules for tables `table1` and `"TaBlE1"` will be the same, and their names can be passed in upper-case or lower-case characters. The same holds for roles. As for databases, before creating a rule, `pg_proaudit` checks whether this database actually exists, and this check is case-sensitive.

The following examples illustrate creation of "weird" object names and specifying correct audit rules for them.

```
create table "TaBlE"'123"();
SELECT pg_proaudit_set_rule(current_database(), null, null, 'public.TabLe"'123',
    null, 'tst cmnt');

create table "taBlE"'"123"();
SELECT pg_proaudit_set_rule(null, null, null, 'public.taBlE"123', null);

create schema "(BIG!) Schema";
```

```
create table "(BIG!) Schema"."My Tab.lE"();
SELECT pg_proaudit_set_rule(null, null, null, '(BIG!) Schema.My Tab.lE', null);

create role "Some '@#$$' person" with login;
SELECT pg_proaudit_set_rule(null, 'authenticate', null, null, 'Some '@#$$' person');

create role "Some '@#$$' person" with superuser;
\c - "Some '@#$$' person"
SELECT pg_proaudit_set_rule(null, 'disconnect', null, null, current_user);

create database " D B 1";
SELECT pg_proaudit_set_rule(' D B 1', 'authenticate', null, null, null);

create database " D b 2";
\c " D b 2"
create extension pg_proaudit;
SELECT pg_proaudit_set_rule(current_database(), 'disconnect', null, null, null);
```

F.6.4. pg_proaudit_settings View

This view displays the current `pg_proaudit` rules, even if they have not been saved into the `pg_proaudit.conf` file yet. The `pg_proaudit_settings` view consists of the following columns:

- `db_name` (text) — name of the database for which to log security events.
- `event_type` (text) — event type to log.
- `object_type` (text) — type of the object for which security events are to be logged.
- `object_name` (text) — name of the object for which security events are to be logged.
- `role_name` (text) — the role on behalf of which logged actions are performed.
- `comment` (text) — comment to describe the created logging rule.

F.6.5. Security Events

You can configure the `pg_proaudit` extension to log classes of security events and specific events by specifying the respective value in the `event_type` argument of the `pg_proaudit_set_rule()` function.

The following classes of security events are supported:

- `ALL_DDL`: CREATE, ALTER, DROP for any database object except stored procedures and functions.
- `ALL_DDL_NONTEMP`: same as `ALL_DDL` but the scope is limited to the objects that are not contained in `pg_temp_nnn` temporary schemas.
- `ALL_DML`: SELECT, INSERT, UPDATE, DELETE, TRUNCATE for any table type; EXECUTE for functions and stored procedures.
- `ALL_DML_NONTEMP`: same as `ALL_DML` but the scope is limited to the objects that are not contained in `pg_temp_nnn` temporary schemas.
- `ALL_MOD`: INSERT, UPDATE, DELETE, TRUNCATE for any table type.
- `ALL_PROC`: CREATE, ALTER, DROP for any function and stored procedure.
- `ALL_ROLE`: CREATE, ALTER, DROP for USER, ROLE, GROUP, PROFILE, as well as execution of the GRANT command.

The following specific security events are supported:

- AUTHENTICATE
- DISCONNECT
- ALTER AGGREGATE
- ALTER COLLATION
- ALTER CONVERSION
- ALTER DATABASE
- ALTER DEFAULT PRIVILEGES
- ALTER DOMAIN
- ALTER EVENT TRIGGER
- ALTER EXTENSION

- ALTER FOREIGN DATA WRAPPER
- ALTER FOREIGN TABLE
- ALTER FUNCTION
- ALTER INDEX
- ALTER LANGUAGE
- ALTER LARGE OBJECT
- ALTER MATERIALIZED VIEW
- ALTER OPERATOR
- ALTER OPERATOR CLASS
- ALTER OPERATOR FAMILY
- ALTER POLICY
- ALTER PROFILE
- ALTER ROLE, ALTER USER, ALTER GROUP
- ALTER RULE
- ALTER SCHEMA
- ALTER SEQUENCE
- ALTER SERVER
- ALTER SYSTEM
- ALTER TABLE
- ALTER TABLESPACE
- ALTER TEXT SEARCH CONFIGURATION
- ALTER TEXT SEARCH DICTIONARY
- ALTER TEXT SEARCH PARSER
- ALTER TEXT SEARCH TEMPLATE
- ALTER TRIGGER
- ALTER TYPE
- ALTER USER MAPPING
- ALTER VIEW
- CLUSTER
- COMMENT
- COPY
- CREATE ACCESS METHOD
- CREATE AGGREGATE
- CREATE CAST
- CREATE COLLATION
- CREATE CONVERSION
- CREATE DATABASE
- CREATE DOMAIN
- CREATE EVENT TRIGGER
- CREATE EXTENSION
- CREATE FOREIGN DATA WRAPPER
- CREATE FOREIGN TABLE
- CREATE FUNCTION
- CREATE INDEX
- CREATE LANGUAGE
- CREATE MATERIALIZED VIEW
- CREATE OPERATOR
- CREATE OPERATOR CLASS
- CREATE OPERATOR FAMILY
- CREATE POLICY
- CREATE PROFILE
- CREATE ROLE, CREATE USER, CREATE GROUP
- CREATE RULE
- CREATE SCHEMA
- CREATE SEQUENCE
- CREATE SERVER

- CREATE TABLE, CREATE TABLE AS, SELECT INTO
- CREATE TABLESPACE
- CREATE TEXT SEARCH CONFIGURATION
- CREATE TEXT SEARCH DICTIONARY
- CREATE TEXT SEARCH PARSER
- CREATE TEXT SEARCH TEMPLATE
- CREATE TRANSFORM
- CREATE TRIGGER
- CREATE TYPE
- CREATE USER MAPPING
- CREATE VIEW
- DEALLOCATE
- DELETE
- DO
- DROP ACCESS METHOD
- DROP AGGREGATE
- DROP CAST
- DROP COLLATION
- DROP CONVERSION
- DROP DATABASE
- DROP DOMAIN
- DROP EVENT TRIGGER
- DROP EXTENSION
- DROP FOREIGN DATA WRAPPER
- DROP FOREIGN TABLE
- DROP FUNCTION
- DROP INDEX
- DROP LANGUAGE
- DROP MATERIALIZED VIEW
- DROP OPERATOR
- DROP OPERATOR CLASS
- DROP OPERATOR FAMILY
- DROP OWNED
- DROP POLICY
- DROP PROFILE
- DROP ROLE, DROP USER, DROP GROUP
- DROP RULE
- DROP SCHEMA
- DROP SEQUENCE
- DROP SERVER
- DROP TABLE
- DROP TABLESPACE
- DROP TEXT SEARCH CONFIGURATION
- DROP TEXT SEARCH DICTIONARY
- DROP TEXT SEARCH PARSER
- DROP TEXT SEARCH TEMPLATE
- DROP TRANSFORM
- DROP TRIGGER
- DROP TYPE
- DROP USER MAPPING
- DROP VIEW
- EXECUTE
- GRANT
- INSERT
- PREPARE
- REASSIGN OWNED

- REFRESH MATERIALIZED VIEW
- REINDEX
- RESET
- REVOKE
- SECURITY LABEL
- SELECT
- SET
- UPDATE
- TRUNCATE TABLE

F.6.6. Security Event Log Configuration Parameters

The `pg_proaudit` extension provides several configuration parameters for managing security event log files. These parameters can be set in the `postgresql.conf` configuration file, or with the help of the `ALTER SYSTEM` command. For the changes to take effect, call the `pg_reload_conf()` function or restart the database server. For additional configuration, the `syslog_ident` and `syslog_facility` configuration parameters can be used.

`pg_proaudit.log_destination (string)`

Defines the method for logging security events. Possible values are:

- `csvlog` — log security events in a CSV file.
- `syslog` — log security events in syslog.

You can specify one or more values separated by commas.

Default: `csvlog`

`pg_proaudit.log_catalog_access (boolean)`

Specifies whether to log access to system catalog objects in the `pg_catalog` schema.

When set to `off`, no events will be logged, even if there are rules for the `CATALOG RELATION`, `CATALOG FUNCTION`, or `ALL` object types.

It is not recommended to set the parameter to `on`, if you do not plan to log events for system catalog objects. Doing otherwise may impact Postgres Pro performance, even if there are no rules for the aforementioned object types.

Default: `off`

`pg_proaudit.log_command_text (boolean)`

Specifies whether to log the SQL command text for security events.

Default: `on`

`pg_proaudit.log_directory (string)`

Specifies the path to the directory that stores CSV log files. This can be an absolute path, or a relative path to the cluster data directory (`PGDATA`). This parameter is used if `pg_proaudit.log_destination` contains the `csvlog` value.

Default: `pg_proaudit`

`pg_proaudit.log_filename (string)`

Defines the filenames of the created security event log files. The filename template can contain %-escapes, similar to the ones listed in the `strftime` specification of the Open Group (<http://pubs.opengroup.org/onlinepubs/009695399/functions/strftime.html>). This parameter is used if `pg_proaudit.log_destination` contains the `csvlog` value.

Default: postgresql-%Y-%m-%d_%H%M%S.log

`pg_proaudit.log_rotation_size` (integer)

Sets the maximum size of the CSV log file, in kilobytes. When this size is achieved, `pg_proaudit` creates a new file for logging security events. This parameter is used if `pg_proaudit.log_destination` contains the `csvlog`. If set to 0, disables size-based creation of new log files.

Default: 10MB

`pg_proaudit.log_rotation_age` (integer)

Sets the maximum lifetime of a log file, in minutes. After this timeframe has elapsed, `pg_proaudit` creates a new file for logging security events. This parameter is used if `pg_proaudit.log_destination` contains the `csvlog` value. If set to 0, disables time-based creation of new log files.

Default: 1 day

`pg_proaudit.log_truncate_on_rotation` (boolean)

Specifies whether to truncate log files when logging is switched to an existing log file. If set to `off`, `pg_proaudit` appends new log entries to the end of the file. This parameter is used if `pg_proaudit.log_destination` contains the `csvlog` value.

Default: `off`

`pg_proaudit.max_rules_count` (integer)

Specifies the maximum number of rules allowed. For the parameter changes to take effect, the database server must be restarted.

Default: 500

F.6.7. Viewing Security Event Log

Security event log files are text files that can be viewed by the operating system facilities. To access log files using SQL, you can use the `file_fdw` extension — a foreign data wrapper for accessing files on the database server. To use this method, complete the following steps:

1. Install the `file_fdw` and create an external server:

```
CREATE EXTENSION file_fdw;  
CREATE SERVER pg_proauditlog FOREIGN DATA WRAPPER file_fdw;
```

2. Create a foreign table, specifying the columns and the absolute path to the log file. The actual log file location is determined by the `pg_proaudit.log_directory` and `pg_proaudit.log_filename` parameters.

```
CREATE FOREIGN TABLE pg_proaudit_log  
    ( log_time timestamp(3) with time zone,  
      current_usr_name text,  
      database_name text,  
      session_pid text,  
      error_severity text,  
      session_line_num bigint,  
      session_line_subcommand_num bigint,  
      event_type text,  
      object_type text,  
      object_name text,  
      status text,  
      error_message text,  
      query_text text,
```

```
query_args text,  
session_usr_name text,  
uuid text,  
xid text,  
vxid text )
```

```
SERVER pg_proauditlog  
OPTIONS (filename 'absolute_file_path_to_log_file.csv', FORMAT 'csv' );
```

Make sure that the `pg_proaudit.log_destination` parameter contains the `csvlog` value, which enables writing security event logs to CSV files.

F.6.8. Examples

As an example, let's set up logging for the following security events:

- authentications/disconnections to the `postgres` database
- all actions of the user if at least one of the user session attributes, `session_user` or `current_user`, is either explicitly set to the `postgres` role, or is directly or indirectly a member of this role
- creating, updating, and deleting any tables
- all operations on the `app_table` table that belongs to the `public` schema

All events must be logged in the CSV format and stored for a week. It is required to set up SQL access to the security event log. To complete the scenario, do the following:

In `psql`, check that the preliminary setup of the `pg_proaudit` extension is complete in the `postgres` database:

```
SHOW shared_preload_libraries;  
shared_preload_libraries  
-----  
pg_proaudit  
  
\dx pg_proaudit  
List of installed extensions  
Name | Version | Schema | Description  
-----+-----+-----+-----  
pg_proaudit | 2.0 | public | provides auditing functionality
```

Add the following lines to the `postgresql.conf` configuration file:

```
pg_proaudit.log_destination = 'csvlog'  
pg_proaudit.log_directory = 'audit'  
pg_proaudit.log_filename = 'audit-%u.csv'  
pg_proaudit.log_rotation_age = 1440  
pg_proaudit.log_rotation_size = 0  
pg_proaudit.log_truncate_on_rotation = on  
pg_proaudit.log_command_text = on
```

For the changes to take effect, run the following query:

```
SELECT pg_reload_conf();
```

Check that the following parameters are set as expected:

```
SHOW pg_proaudit.log_destination;  
SHOW pg_proaudit.log_directory;  
SHOW pg_proaudit.log_filename;  
SHOW pg_proaudit.log_rotation_age;  
SHOW pg_proaudit.log_rotation_size;  
SHOW pg_proaudit.log_truncate_on_rotation;  
SHOW pg_proaudit.log_command_text;
```


Suppose your PGDATA environment variable points to the cluster data directory. Since the `pg_proaudit.log_directory` defines a relative path to the log files, they will be located in the `$PGDATA/audit` directory. Let's create an empty file for each day of the week and make them available to their owner only:

```
touch $PGDATA/audit/audit-1.csv
touch $PGDATA/audit/audit-2.csv
touch $PGDATA/audit/audit-3.csv
touch $PGDATA/audit/audit-4.csv
touch $PGDATA/audit/audit-5.csv
touch $PGDATA/audit/audit-6.csv
touch $PGDATA/audit/audit-7.csv
chmod 600 $PGDATA/audit/audit-*.csv
```

Create a table for reading log entries:

```
CREATE TABLE pg_proaudit_log (
    log_time timestamp(3) with time zone,
    current_usr_name text,
    database_name text,
    session_pid text,
    error_severity text,
    session_line_num bigint,
    session_line_subcommand_num bigint,
    event_type text,
    object_type text,
    object_name text,
    status text,
    error_message text,
    query_text text,
    query_args text,
    session_usr_name text,
    uuid text,
    xid text,
    vxid text
);
```

Install the `file_fdw` extension and create an external server:

```
CREATE EXTENSION file_fdw;
CREATE SERVER pg_proauditlog FOREIGN DATA WRAPPER file_fdw;
```

Now let's create seven child foreign tables for the `pg_proaudit_log` table, for each day of the week:

```
CREATE FOREIGN TABLE pg_proaudit_log_1 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
    OPTIONS (filename '/path_to_PGDATA/audit/audit-1.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_2 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
    OPTIONS (filename '/path_to_PGDATA/audit/audit-2.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_3 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
    OPTIONS (filename '/path_to_PGDATA/audit/audit-3.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_4 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
    OPTIONS (filename '/path_to_PGDATA/audit/audit-4.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_5 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
    OPTIONS (filename '/path_to_PGDATA/audit/audit-5.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_6 () INHERITS (pg_proaudit_log) SERVER
pg_proauditlog
```

Additional Supplied Modules and Extensions Shipped in postgrespro-sdm-17-contrib

```

OPTIONS (filename '/path_to_PGDATA/audit/audit-6.csv', FORMAT 'csv');
CREATE FOREIGN TABLE pg_proaudit_log_7 () INHERITS (pg_proaudit_log) SERVER
pg_proaudit_log
OPTIONS (filename '/path_to_PGDATA/audit/audit-7.csv', FORMAT 'csv');

```

To set up logging for the required security events, connect to the postgres database and execute the following commands:

```

SELECT pg_proaudit_set_rule (current_database(), 'AUTHENTICATE', null, null, null, 'Any
authentication in the current DB');
SELECT pg_proaudit_set_rule (current_database(), 'DISCONNECT', null, null, null, 'Any
disconnect from the current DB');
SELECT pg_proaudit_set_rule (current_database(), 'ALL', 'TABLE', null, null, 'Any
operations with any table in the current DB');
SELECT pg_proaudit_set_rule (current_database(), 'ALL', null, null, 'postgres', 'Any
operation by "postgres" user in the current DB');

```

Create the app_table table and enable logging for all operations on this table:

```

CREATE TABLE app_table (id int, name text);
SELECT pg_proaudit_set_rule (current_database(), 'ALL', null, 'public.app_table',
null);

```

Check that event logging is configured as expected:

```

SELECT * FROM pg_proaudit_settings;

```

db_name	event_type	object_type	object_name	role_name	comment
postgres	authenticate	ALL			Any authentication in the current DB
postgres	disconnect	ALL			Any disconnect from the current DB
postgres	ALL	table			Any operations with any table in the current DB
postgres	ALL	ALL		postgres	Any operation by "postgres" user in the current DB
postgres	ALL	ALL	public.app_table		

(5 rows)

Save these logging rules into the pg_proaudit.conf file, so that they are not lost after the server restart:

```

SELECT pg_proaudit_save();

```

Let's run several queries on the app_table table:

```

INSERT INTO app_table VALUES (1, 'first');
SELECT * FROM app_table;

```

Check the log entries for the app_table table:

```

SELECT to_char(log_time, 'DD.MM.YY HH24:MI:SS') AS when, current_usr_name,
session_pid, event_type, query_text, session_usr_name
FROM pg_proaudit_log
WHERE object_name = 'public.app_table';

```

when	current_usr_name	session_pid	event_type	query_text	session_usr_name
27.09.23 12:44:27	postgres	2010	CREATE TABLE	CREATE TABLE app_table (id int, name text);	postgres

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

27.09.23 12:45:55	postgres	2010	INSERT	INSERT INTO
app_table VALUES (1, 'first'); postgres				
27.09.23 12:46:00	postgres	2010	SELECT	SELECT * FROM
app_table; postgres				
(3 rows)				

We have set up a weekly rotation of log files, with the log file switched once a day. It means that queries to the `pg_proaudit_log` table will return only those security events that have happened in the latest week. Older events will be automatically removed at log file rotation. To define additional access constraints for specific log entries, you can create separate views based on queries to the `pg_proaudit_log` table and grant read rights to such views using built-in Postgres Pro access control methods.

F.7. pgpro_stats — a means for tracking planning and execution statistics of all SQL statements executed by a server

The `pgpro_stats` extension provides a means for tracking planning and execution statistics of all SQL statements executed by a server. It is based on the `pg_stat_statements` module and provides the following additional functionality:

- Storing query plans in addition to query statements.
- Configuring sample rate for statistics collection to reduce overhead.
- Calculating wait event statistics for executed queries.
- Calculating resource usage statistics of statement planning and execution.
- Calculating cache invalidation statistics.
- Calculating additional archiver statistics.
- Providing an interface to statistics about vacuuming databases, tables, and indexes collected by the core system.
- Tracing of application sessions.
- Creating views that emulate other extensions.

Note

When the server gets shut down, `pgpro_stats` saves the collected statistics to two dump files on disk in the format that depends on the number and order of columns in the `pgpro_stats_statements` and `pgpro_stats_totals` views. If this format has changed in a new `pgpro_stats` version, during the first server restart after upgrading the extension, the saved statistics will fail to be read with an error reported to the log:

```
LOG:  pgpro_stats: could not load statements data file of obsolete format
      "pg_stat/pgpro_stats.stat"
```

or

```
LOG:  pgpro_stats: could not load statements data file of obsolete format
      "pg_stat/pgpro_stats_totals.stat"
```

The contents of the files will be ignored, and respective statistics will be nullified. To learn whether the format has changed, see the [Release Notes](#) for `pgpro_stats` version being installed. To retain the previously collected statistics, download them as CSV files before upgrading.

The background information, along with views and types, related to calculating cache invalidation statistics is provided in a separate section [Section F.7.9](#).

F.7.1. Limitations

- `pgpro_stats` can sometimes fail to match identical parameters in the query statement and the corresponding query plan.
- Some SPI queries are not included into statistics.
- Texts and plans of some SPI queries are not normalized.
- `pgpro_stats` is incompatible with `pg_stat_statements`, as well as other extensions that use parser, planner, or executor hooks to modify parse and plan trees and execution of the queries. Moreover, if both `pgpro_stats` and `pg_stat_statements` are on the list of `shared_preload_libraries`, the database server will not start. Note also that in order to dump the final versions of the queries and plans, `pgpro_stats` should be the last on the list of `shared_preload_libraries`, but some existing extensions may not work at all unless they are the last on this list.

- `pgpro_stats` may not work correctly with third-party extensions that produce `CustomScan` and `ForeignScan` nodes.

F.7.2. Installation and Setup

`pgpro_stats` is provided with Postgres Pro Shardman as a separate pre-built package `pgpro-stats-sdm-17` (for the detailed installation instructions, see [Chapter 17](#)). Once you have `pgpro_stats` installed, complete the following steps to enable `pgpro_stats`:

1. Add `pgpro_stats` to the [shared_preload_libraries](#) parameter in the `postgresql.conf` file:

Important

Ensure that in `shared_preload_libraries`, `pgpro_stats` is listed after `shardman`. For example:

```
shared_preload_libraries = 'postgres_fdw, shardman, pgpro_stats'
```

2. Restart the Postgres Pro Shardman instance for the changes to take effect.

Once the server is reloaded, `pgpro_stats` starts tracking statistics across all databases of the cluster. If required, you can change the scope of statistics collection or disable it altogether using `pgpro_stats` [configuration parameters](#).

3. To access the collected statistics, you have to create `pgpro_stats` extension:

```
CREATE EXTENSION pgpro_stats;
```

In addition, query identifier calculation must be enabled in order for `pgpro_stats` to be active, which is done automatically if [compute_query_id](#) is set to `auto` or `on`, or any third-party module that calculates query identifiers is loaded.

F.7.3. Usage

- [Collecting Statistics on Query Statements and Plans](#)
- [Collecting Distributed Statement Statistics in Postgres Pro Shardman](#)
- [Monitoring Custom Metrics](#)

F.7.3.1. Collecting Statistics on Query Statements and Plans

Once installed, the `pgpro_stats` extension starts collecting statistics on the executed statements. The collected data is similar to the one provided by [pg_stat_statements](#), but also includes information on query plans and wait events for each query type. The statistics is saved into an in-memory ring buffer and is accessible through the [pgpro_stats_statements](#) view.

By default, `pgpro_stats` collects statistics on all the executed statements that satisfy the [pgpro_stats.track](#) and [pgpro_stats.track_utility](#) settings. If performance is a concern, you can set a sample rate for queries using the [pgpro_stats.query_sample_rate](#) parameter, and `pgpro_stats` will randomly select queries for statistics calculation at the specified rate.

To collect statistics on wait events, `pgpro_stats` uses time-based sampling. Wait events are sampled at the time interval specified by the [pgpro_stats.profile_period](#) parameter, which is set to 10ms by default. If the sampling shows that the process is waiting, the [pgpro_stats.profile_period](#) value is added to the wait event duration. Thus, time estimation for each wait event remains valid even if the [pgpro_stats.profile_period](#) parameter value has changed. If you are not interested in wait event statistics, you can disable wait event sampling by setting the [pgpro_stats.enable_profile](#) parameter to `false`.

`pgpro_stats_statements.plans` and `pgpro_stats_statements.calls` aren't always expected to match because planning and execution statistics are updated at their respective end phase, and only for successful operations. For example, if a statement is successfully planned but fails during the execution

phase, only its planning statistics will be updated. If planning is skipped because a cached plan is used, only its execution statistics will be updated.

As an example, let's create a table with some random data and build an index on this table:

```
CREATE TABLE test AS (SELECT i, random() x FROM generate_series(1,1000000) i);
CREATE INDEX test_x_idx ON test (x);
```

Now run the following query several times using different values for `:x_min` and `:x_max`:

```
select * from test where x >= :x_min and x <= :x_max;
```

The collected statistics should appear in the `pgpro_stats_statements` view:

```
SELECT queryid, query, planid, plan, wait_stats FROM pgpro_stats_statements WHERE query
LIKE 'select * from test where%';
```

```
-[ RECORD
```

```
1 ]-----
queryid      | 1109491335754870054
query        | select * from test where x >= $1 and x <= $2
planid       | 8287793242828473388
plan         | Gather
              |   Output: i, x
              |   Workers Planned: 2
              |   -> Parallel Seq Scan on public.test
              |       Output: i, x
              |       Filter: ((test.x >= $3) AND (test.x <= $4))
wait_stats   | {"IO": {"DataFileRead": 10}, "IPC": {"BgWorkerShutdown": 10}, "Total":
              | {"IO": 10, "IPC": 10, "Total": 20}}
```

```
-[ RECORD
```

```
2 ]-----
queryid      | 1109491335754870054
query        | select * from test where x >= $1 and x <= $2
planid       | -9045072158333552619
plan         | Bitmap Heap Scan on public.test
              |   Output: i, x
              |   Recheck Cond: ((test.x >= $3) AND (test.x <= $4))
              |   -> Bitmap Index Scan on test_x_idx
              |       Index Cond: ((test.x >= $5) AND (test.x <= $6))
wait_stats   | {"IO": {"DataFileRead": 40}, "Total": {"IO": 40, "Total": 40}}
```

```
-[ RECORD
```

```
3 ]-----
queryid      | 1109491335754870054
query        | select * from test where x >= $1 and x <= $2
planid       | -1062789671372193287
plan         | Seq Scan on public.test
              |   Output: i, x
              |   Filter: ((test.x >= $3) AND (test.x <= $4))
wait_stats   | NULL
```

```
-[ RECORD
```

```
4 ]-----
queryid      | 1109491335754870054
query        | select * from test where x >= $1 and x <= $2
planid       | -1748292253893834280
plan         | Index Scan using test_x_idx on public.test
              |   Output: i, x
              |   Index Cond: ((test.x >= $3) AND (test.x <= $4))
```

```
|  
wait_stats | NULL
```

F.7.3.2. Collecting Distributed Statement Statistics in Postgres Pro Shardman

In Postgres Pro Shardman, distributed statement statistics are collected as explained below.

During execution of distributed queries, Postgres Pro Shardman sends derived SQL queries to remote nodes that hold data partitions involved in the query execution. Let's call these SQL queries query fragments. Postgres Pro Shardman sends such queries using the [postgres_fdw](#) extension. The node that queries the sharded table is called the *coordinator*, while in this context, the nodes that accept query fragments are called *shard nodes*.

When the `pgpro_stats` extension is enabled on a Postgres Pro Shardman cluster node, it collects statistics about local and distributed queries. The information about distributed queries initiated by this node is incomplete because it misses data about remote query fragments. The statistics concerning queries initiated by other nodes is also ambiguous because there is no simple way for a user to determine the distributed query to which the fragment corresponds.

To address these issues, `pgpro_stats` introduces an aggregation of statistics for the distributed queries. These aggregated statistics can be accessed with the `pgpro_stats_sdm_statements` view. However, each Postgres Pro Shardman node collects statistics for all the statements, so that the `pgpro_stats_statements` view can work the way it did before.

When a node receives a query fragment, it saves its statistics to a separate shared hash table. Periodically and asynchronously, each node sends this information from a separate table to the coordinator corresponding to the query. The coordinator aggregates the statistical data obtained from the query fragments with the statistics of its parent query, which is the query initiated by the client.

The `pgpro_stats` extension starts a separate background worker. This worker is responsible for sending the accumulated statistics to the coordinator nodes periodically or when triggered by the guard latch. The collecting function sets this latch when the hash table is almost full.

To reduce the network traffic initiated by a statistics sender, compression is applied to the statistics data sent. The compression method can be selected by the `pgpro_stats.transport_compression` configuration parameter.

Each node stores the total number of statistics entries received from the shard node and the timestamp of when they were last received. When a coordinator node receives a statistics message, it updates the appropriate values, which are accessible using the SQL interface by the `pgpro_stats_sdm_stats_updated` function.

F.7.3.3. Monitoring Custom Metrics

With `pgpro_stats`, you can define custom metrics to be monitored. The collected data will be saved into an in-memory ring buffer and then sent to a monitoring system. Unlike direct polling of a database by a monitoring system that can lose some data if the connection is interrupted, this approach allows to get all the collected data regardless of connection issues, as long as this data is still available in the ring buffer.

To set up a custom metric to collect, do the following:

1. For each metric, define all configuration parameters listed in [Section F.7.7.2](#). You must specify a unique numeric identifier of each metric in the parameter names.

For example, to monitor index bloating each 60 seconds, you can define a new metric by setting metrics-related parameters as follows:

```
pgpro_stats.metric_1_name = index_bloat  
pgpro_stats.metric_1_query = 'select iname, ibloat, ipages from bloat'  
pgpro_stats.metric_1_db = 'postgres'  
pgpro_stats.metric_1_user = postgres  
pgpro_stats.metric_1_period = '60s'
```

2. Restart the server.

`pgpro_stats` starts collecting statistics on executed statements and saves it into the ring buffer, and the collected data appears in the `pgpro_stats_metrics` view:

```
SELECT * FROM pgpro_stats_metrics;
```

Once the new metric is added, its parameters can be changed without a server restart by simply reloading the `postgresql.conf` configuration file.

3. If required, set up data export to a monitoring system of your choice.

F.7.4. Views

F.7.4.1. The `pgpro_stats_statements` View

The statistics gathered by the module are made available via a view named `pgpro_stats_statements`. This view contains one row for each distinct database ID, user ID and query ID (up to the maximum number of distinct statements that the module can track). The columns of the view are shown in [Table F.8](#).

Table F.8. `pgpro_stats_statements` Columns

Name	Type	References	Description
<code>userid</code>	<code>oid</code>	<code>pg_authid .oid</code>	OID of user who executed the statement
<code>dbid</code>	<code>oid</code>	<code>pg_database .oid</code>	OID of database in which the statement was executed
<code>toplevel</code>	<code>bool</code>		True if the query was executed as a top-level statement (always true if <code>pgpro_stats.track</code> is set to <code>top</code>)
<code>queryid</code>	<code>bigint</code>		Internal hash code, computed from the statement's parse tree
<code>planid</code>	<code>bigint</code>		Internal hash code, computed from the statement's plan tree
<code>query</code>	<code>text</code>		Text of a representative statement
<code>plan</code>	<code>text</code>		The text of the query plan, in the format defined by the <code>pgpro_stats.plan_format</code> configuration parameter
<code>plans</code>	<code>int8</code>		Number of times the statement was planned (if <code>pgpro_stats.track_planning</code> is enabled, otherwise zero)
<code>total_plan_time</code>	<code>float8</code>		Total time spent planning the statement, in milliseconds (if <code>pgpro_stats.track_planning</code> is enabled, otherwise zero).

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
min_plan_time	float8		Minimum time spent planning the statement, in milliseconds (if <code>pg-pro_stats.track_planning</code> is enabled, otherwise zero)
max_plan_time	float8		Maximum time spent planning the statement, in milliseconds (if <code>pg-pro_stats.track_planning</code> is enabled, otherwise zero)
mean_plan_time	float8		Mean time spent planning the statement, in milliseconds (if <code>pgpro_stats.track_planning</code> is enabled, otherwise zero)
stddev_plan_time	float8		Population standard deviation of time spent planning the statement, in milliseconds (if <code>pg-pro_stats.track_planning</code> is enabled, otherwise zero)
plan_rusage	pgpro_stats_rusage		Resource usage statistics of the statement planning.
calls	int8		Number of times the statement was executed
total_exec_time	float8		Total time spent executing the statement, in milliseconds
min_exec_time	float8		Minimum time spent executing the statement, in milliseconds
max_exec_time	float8		Maximum time spent executing the statement, in milliseconds
mean_exec_time	float8		Mean time spent executing the statement, in milliseconds
stddev_exec_time	float8		Population standard deviation of time spent executing the statement, in milliseconds
exec_rusage	pgpro_stats_rusage		Resource usage statistics of the statement execution.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
rows	int8		Total number of rows retrieved or affected by the statement
shared_blks_hit	int8		Total number of shared block cache hits by the statement
shared_blks_read	int8		Total number of shared blocks read by the statement
shared_blks_dirtied	int8		Total number of shared blocks dirtied by the statement
shared_blks_written	int8		Total number of shared blocks written by the statement
local_blks_hit	int8		Total number of local block cache hits by the statement
local_blks_read	int8		Total number of local blocks read by the statement
local_blks_dirtied	int8		Total number of local blocks dirtied by the statement
local_blks_written	int8		Total number of local blocks written by the statement
temp_blks_read	int8		Total number of temp blocks read by the statement
temp_blks_written	int8		Total number of temp blocks written by the statement
shared_blk_read_time	float8		Total time the statement spent reading shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
shared_blk_write_time	float8		Total time the statement spent writing shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
local_blk_read_time	float8		Total time the statement spent reading local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
local_blk_write_time	float8		Total time the statement spent writing local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.
temp_blk_read_time	float8		Total time the statement spent reading temp blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.
temp_blk_write_time	float8		Total time the statement spent writing temp blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.
wal_records	int8		Total number of WAL records generated by the statement
wal_fpi	int8		Total number of WAL full page images generated by the statement
wal_bytes	numeric		Total amount of WAL bytes generated by the statement
jit_functions	int8		Total number of functions JIT-compiled by the statement. In Postgres Pro versions lower than 15, contains zero.
jit_generation_time	float8		Total time spent by the statement on generating JIT code, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_inlining_count	int8		Number of times functions used in the statement have been inlined. In Postgres Pro versions lower than 15, contains zero.
jit_inlining_time	float8		Total time spent by the statement on inlining functions, in millisec-

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
			onds. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_count	int8		Number of times the statement has been optimized. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_time	float8		Total time spent by the statement on optimizing, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_emission_count	int8		Number of times code has been emitted by the statement. In Postgres Pro versions lower than 15, contains zero.
jit_emission_time	float8		Total time spent by the statement on emitting code, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_deform_count	int8		Total number of tuple deform functions JIT-compiled by the statement. In Postgres Pro versions lower than 17, contains zero.
jit_deform_time	float8		Total time spent by the statement on JIT-compiling tuple deform functions, in milliseconds. In Postgres Pro versions lower than 17, contains zero.
wait_stats	jsonb		A <code>jsonb</code> object containing statistics on wait events, for each execution of the query that uses the corresponding plan. Each statistic is provided in milliseconds and is a multiple of the pgpro_stats.profile_period configuration parameter.
inval_msgs	pgpro_stats_inval_msgs		Number of cache invalidation messages by type generated by the statement (if this is support-

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
			ed by the server, otherwise zero).
stats_since	timestamp with time zone		Time at which statistics gathering started for this statement
minmax_stats_since	timestamp with time zone		Time at which min/max statistics gathering started for this statement (fields min_plan_time , max_plan_time , min_exec_time and max_exec_time)

Take into account that like [pg_stat_statements](#), pgpro_stats normalizes into one record those DML queries (containing SELECT, INSERT, UPDATE, DELETE and MERGE commands) that have equivalent structures according to some internal hash value. Being compared this way, two queries are normally considered equal if they are semantically equivalent up to constants included in the queries. All the other commands are, however, compared strictly as query texts. When the value of a constant in a query is ignored for comparison with other queries, this constant is replaced in the pgpro_stats output with a symbol of a parameter, such as, \$k, where k is a positive integer. If a query already contains parameters, the initial value of k equals the number following the last number of a \$n parameter in the original query text. If there are no parameters, the initial value of k equals 1. Note that sometimes hidden parameter symbols affect this numbering. For example, PL/pgSQL uses such hidden symbols to insert values of function local variables into queries, so a PL/pgSQL statement SELECT i + 1 INTO j will be represented as SELECT i + \$2 in the normalized query text.

pgpro_stats uses a similar technique to normalize plan texts. When doing so, an attempt is made to associate numbers of constants in the plan text with the corresponding numbers of constants in the query text. If such an attempt appears unsuccessful for a certain constant in the plan text, it is assigned the number following the maximum number of a constant replaced in the query text. For example, consider the query:

```
SELECT 1::int, 'abc'::VARCHAR(3), 2::int;
```

pgpro_stats will replace numbers of constants in the query text and in the text of the corresponding plan as follows:

```
postgres=# SELECT query, plan FROM pgpro_stats_statements;
               query                               |                               plan
-----
+-----+
SELECT $1::int, $2::VARCHAR(3), $3::int          | Result
      +                                           |
                                           |   Output: $1, $4, $3
      +                                           |
```

In this plan text, it appeared possible to associate constants numbered 1 and 3 from the query text, but not the constant numbered 2, and the latter was replaced with the number following the maximum number in the query text, that is, number 4.

Replacement of numbers in plan texts has an exception for version numbers of XML documents. If in the original query such a number is represented with a constant, e.g., '1.0', it is retained as is in the plan text rather than replaced with \$k. If the version number of an XML document is represented with an expression, replacement of constants follows usual rules.

F.7.4.2. The pgpro_stats_sdm_statements View

Unlike [pgpro_stats_statements](#), the pgpro_stats_sdm_statements view displays statistics only for distributed queries executed from a coordinator. These statistics are presented in an aggregated form as explained in [Section F.7.3.2](#). In addition to pgpro_stats_statements, this view also has several columns

with data related to execution time and network traffic generated as a result of cluster interactions. The columns of the view are shown in [Table F.9](#).

Table F.9. pgpro_stats_sdm_statements Columns

Name	Type	References	Description
userid	oid	pg_authid .oid	OID of user who executed the statement
dbid	oid	pg_database .oid	OID of database in which the statement was executed
toplevel	bool		True if the query was executed as a top-level statement (always true if pgpro_stats.track is set to top)
queryid	bigint		Internal hash code, computed from the statement's parse tree
planid	bigint		Internal hash code, computed from the statement's plan tree
query	text		Text of a representative statement
plan	text		The text of the query plan, in the format defined by the pgpro_stats.plan_format configuration parameter
plans	int8		Number of times the statement was planned (if pgpro_stats.track_planning is enabled, otherwise zero)
total_plan_time	float8		Total time spent planning the statement, in milliseconds (if pgpro_stats.track_planning is enabled, otherwise zero).
min_plan_time	float8		Minimum time spent planning the statement, in milliseconds (if pgpro_stats.track_planning is enabled, otherwise zero)
max_plan_time	float8		Maximum time spent planning the statement, in milliseconds (if pgpro_stats.track_planning is enabled, otherwise zero)
mean_plan_time	float8		Mean time spent planning the statement, in

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
			milliseconds (if <code>pgpro_stats.track_planning</code> is enabled, otherwise zero)
<code>stddev_plan_time</code>	<code>float8</code>		Population standard deviation of time spent planning the statement, in milliseconds (if <code>pgpro_stats.track_planning</code> is enabled, otherwise zero)
<code>plan_rusage</code>	pgpro_stats_rusage		Resource usage statistics of the statement planning.
<code>calls</code>	<code>int8</code>		Number of times the statement was executed
<code>total_exec_time</code>	<code>float8</code>		Total time spent executing the statement, in milliseconds
<code>min_exec_time</code>	<code>float8</code>		Minimum time spent executing the statement, in milliseconds
<code>max_exec_time</code>	<code>float8</code>		Maximum time spent executing the statement, in milliseconds
<code>mean_exec_time</code>	<code>float8</code>		Mean time spent executing the statement, in milliseconds
<code>stddev_exec_time</code>	<code>float8</code>		Population standard deviation of time spent executing the statement, in milliseconds
<code>exec_rusage</code>	pgpro_stats_rusage		Resource usage statistics of the statement execution.
<code>rows</code>	<code>int8</code>		Total number of rows retrieved or affected by the statement
<code>shared_blks_hit</code>	<code>int8</code>		Total number of shared block cache hits by the statement
<code>shared_blks_read</code>	<code>int8</code>		Total number of shared blocks read by the statement
<code>shared_blks_dirtied</code>	<code>int8</code>		Total number of shared blocks dirtied by the statement
<code>shared_blks_written</code>	<code>int8</code>		Total number of shared blocks written by the statement

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
local_blks_hit	int8		Total number of local block cache hits by the statement
local_blks_read	int8		Total number of local blocks read by the statement
local_blks_dirtied	int8		Total number of local blocks dirtied by the statement
local_blks_written	int8		Total number of local blocks written by the statement
temp_blks_read	int8		Total number of temp blocks read by the statement
temp_blks_written	int8		Total number of temp blocks written by the statement
shared_blk_read_time	float8		Total time the statement spent reading shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
shared_blk_write_time	float8		Total time the statement spent writing shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
local_blk_read_time	float8		Total time the statement spent reading local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.
local_blk_write_time	float8		Total time the statement spent writing local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.
temp_blk_read_time	float8		Total time the statement spent reading temp blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
temp_blk_write_time	float8		Total time the statement spent writing temp blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.
wal_records	int8		Total number of WAL records generated by the statement
wal_fpi	int8		Total number of WAL full page images generated by the statement
wal_bytes	numeric		Total amount of WAL bytes generated by the statement
fdw_recv_bytes	numeric		Size of the incoming network traffic for intra-cluster data exchange, in bytes
fdw_sent_bytes	numeric		Size of the outgoing network traffic for intra-cluster data exchange, in bytes
client_recv_bytes	numeric		Size of the incoming network traffic for data exchange with a client application, in bytes. Deprecated and contains zero.
client_sent_bytes	numeric		Size of the outgoing network traffic for data exchange with a client application, in bytes. Deprecated and contains zero.
fdw_wait_time	numeric		Wait time for a response from other cluster nodes during processing of this query, in milliseconds
jit_functions	int8		Total number of functions JIT-compiled by the statement. In Postgres Pro versions lower than 15, contains zero.
jit_generation_time	float8		Total time spent by the statement on generating JIT code, in milliseconds. In Postgres Pro

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
			versions lower than 15, contains zero.
jit_inlining_count	int8		Number of times functions used in the statement have been inlined. In Postgres Pro versions lower than 15, contains zero.
jit_inlining_time	float8		Total time spent by the statement on inlining functions, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_count	int8		Number of times the statement has been optimized. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_time	float8		Total time spent by the statement on optimizing, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_emission_count	int8		Number of times code has been emitted by the statement. In Postgres Pro versions lower than 15, contains zero.
jit_emission_time	float8		Total time spent by the statement on emitting code, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_deform_count	int8		Total number of tuple deform functions JIT-compiled by the statement. In Postgres Pro versions lower than 17, contains zero.
jit_deform_time	float8		Total time spent by the statement on JIT-compiling tuple deform functions, in milliseconds. In Postgres Pro versions lower than 17, contains zero.
wait_stats	jsonb		A <code>jsonb</code> object containing statistics on wait events, for each execution of the query that uses the corresponding

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	References	Description
			plan. Each statistic is provided in milliseconds and is a multiple of the pgpro_stats.profile_period configuration parameter.
inval_msgs	pgpro_stats_inval_msgs		Number of cache invalidation messages by type generated by the statement (if this is supported by the server, otherwise zero).
stats_since	timestamp with time zone		Time at which statistics gathering started for this statement
minmax_stats_since	timestamp with time zone		Time at which min/max statistics gathering started for this statement (fields <code>min_plan_time</code> , <code>max_plan_time</code> , <code>min_exec_time</code> and <code>max_exec_time</code>)

F.7.4.3. The `pgpro_stats_totals` View

The aggregate statistics gathered by the module are made available via a view named `pgpro_stats_totals`. This view contains one row for each distinct object (up to the maximum number of distinct objects that the module can track). The columns of the view are shown in [Table F.10](#).

Table F.10. `pgpro_stats_totals` Columns

Name	Type	Description
object_type	text	Type of the object for which aggregated statistics are collected: "cluster", "database", "user", "client_addr", "application", "backend", "session"
object_id	bigint	ID of the object: oid for databases and users, pid for backends, sid for sessions, NULL for others
object_name	text	Textual name of the object or NULL
queries_planned	int8	Number of queries planned
total_plan_time	float8	Total time spent in the planning of statements, in milliseconds
total_plan_rusage	pgpro_stats_rusage	Aggregate resource usage statistics of the statement planning
queries_executed	int8	Number of queries executed
total_exec_time	float8	Total time spent in the execution of statements, in milliseconds
total_exec_rusage	pgpro_stats_rusage	Aggregate resource usage statistics of the statement execution

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
rows	int8	Total number of rows retrieved or affected by the statements
shared_blks_hit	int8	Total number of shared block cache hits by the statements
shared_blks_read	int8	Total number of shared blocks read by the statements
shared_blks_dirtied	int8	Total number of shared blocks dirtied by the statements
shared_blks_written	int8	Total number of shared blocks written by the statements
local_blks_hit	int8	Total number of local block cache hits by the statements
local_blks_read	int8	Total number of local blocks read by the statements
local_blks_dirtied	int8	Total number of local blocks dirtied by the statements
local_blks_written	int8	Total number of local blocks written by the statements
temp_blks_read	int8	Total number of temp blocks read by the statements
temp_blks_written	int8	Total number of temp blocks written by the statements
shared_blk_read_time	float8	Total time the statements spent reading shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
shared_blk_write_time	float8	Total time the statements spent writing shared blocks, in milliseconds (if track_io_timing is enabled, otherwise zero)
local_blk_read_time	float8	Total time the statements spent reading local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.
local_blk_write_time	float8	Total time the statements spent writing local blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 17, contains zero.
temp_blk_read_time	float8	Total time the statements spent reading temp blocks, in milliseconds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.
temp_blk_write_time	float8	Total time the statements spent writing temp blocks, in milliseconds

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
		onds (if track_io_timing is enabled, otherwise zero). In Postgres Pro versions lower than 15, contains zero.
wal_records	int8	Total number of WAL records generated by the statements
wal_fpi	int8	Total number of WAL full page images generated by the statements
wal_bytes	numeric	Total amount of WAL bytes generated by the statements
jit_functions	int8	Total number of functions JIT-compiled by the statements. In Postgres Pro versions lower than 15, contains zero.
jit_generation_time	float8	Total time spent by the statements on generating JIT code, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_inlining_count	int8	Number of times functions used in the statements have been inlined. In Postgres Pro versions lower than 15, contains zero.
jit_inlining_time	float8	Total time spent by the statements on inlining functions, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_count	int8	Number of times the statements have been optimized. In Postgres Pro versions lower than 15, contains zero.
jit_optimization_time	float8	Total time spent by the statements on optimizing, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_emission_count	int8	Number of times code has been emitted by the statements. In Postgres Pro versions lower than 15, contains zero.
jit_emission_time	float8	Total time spent by the statements on emitting code, in milliseconds. In Postgres Pro versions lower than 15, contains zero.
jit_deform_count	int8	Total number of tuple deform functions JIT-compiled by the statements. In Postgres Pro versions lower than 17, contains zero.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
jit_deform_time	float8	Total time spent by the statements on JIT-compiling tuple deform functions, in milliseconds. In Postgres Pro versions lower than 17, contains zero.
wait_stats	jsonb	A jsonb object containing statistics on wait events for each execution of the queries. Each statistic is provided in milliseconds and is a multiple of the pgpro_stats.profile_period configuration parameter.
inval_msgs	pgpro_stats_inval_msgs	Number of cache invalidation messages by type generated by the statements (if this is supported by the server, otherwise zero).
cache_resets	int4	Number of shared cache resets (only for cluster, databases and backends). Gets incremented for a backend when it receives a full cache reset message.
stats_since	timestamp with time zone	Time at which statistics gathering started for the statements

F.7.4.4. The pgpro_stats_info View

The statistics of the `pgpro_stats` module itself are tracked and made available via a view named `pgpro_stats_info`. This view contains only a single row. The columns of the view are shown in [Table F.11](#).

Table F.11. pgpro_stats_info Columns

Name	Type	Description
dealloc	bigint	Total number of times <code>pgpro_stats_statements</code> entries about the least-executed statements were deallocated because more distinct statements than pgpro_stats.max were observed
stats_reset	timestamp with time zone	Time at which all statistics in the pgpro_stats_statements view were last reset

F.7.4.5. The pgpro_stats_metrics View

The metrics gathered by `pgpro_stats` are displayed in the `pgpro_stats_metrics` view. The table below describes the columns of the view.

Table F.12. pgpro_stats_metrics Columns

Name	Type	Description
metric_number	int4	A unique ID of the collected metric assigned by user. This ID is included into parameter names that define the metric.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
metric_name	text	The name of the metric defined by the <code>pgpro_stats.metric_N_name</code> parameter
db_name	text	The name of the database for which a particular metric was collected
ts	timestampz	The time when the metric value got calculated
value	jsonb	The result of the query used for metric measurement. It is serialized in <code>jsonb</code> as an array of objects received via <code>to_jsonb(resulting_row)</code> . If an error occurs, a single object is returned that contains <code>code</code> , <code>message</code> , <code>detail</code> , and <code>hint</code> fields.

F.7.4.6. The `pgpro_stats_archiver` View

The `pgpro_stats_archiver` view will contain one row showing data about the archiver process of the cluster.

Table F.13. `pgpro_stats_archiver` Columns

Column	Type	Description
archived_count	bigint	Number of WAL files that have been successfully archived
last_archived_wal	text	Name of the last WAL file successfully archived
last_archived_time	timestamp with time zone	Time of the last successful archive operation
failed_count	bigint	Number of failed attempts for archiving WAL files
last_failed_wal	text	Name of the WAL file of the last failed archival operation
last_failed_time	timestamp with time zone	Time of the last failed archival operation
active_time	int8	Overall time that the archiver process was active
archive_command_time	int8	Overall execution time of the archive command
stats_reset	timestamp with time zone	Time at which these statistics were last reset

F.7.4.7. The `pgpro_stats_vacuum_database` View

Important

Starting with Postgres Pro 16, this view contains no data because the statistics to be displayed are available through the catalog view `pg_stats_vacuum_database` (see [System Views](#) for details).

The `pgpro_stats_vacuum_database` view will contain one row for each database in the current cluster, showing statistics about vacuuming that database. These statistics are collected by the core system as explained in [Section 27.2](#). The table below describes the columns of the view.

Table F.14. `pgpro_stats_vacuum_database` Columns

Column	Type	Description
<code>dbid</code>	<code>oid</code>	OID of a database
<code>total_blks_read</code>	<code>int8</code>	Number of database blocks read by vacuum operations performed on this database
<code>total_blks_hit</code>	<code>int8</code>	Number of times database blocks were found in the buffer cache by vacuum operations performed on this database
<code>total_blks_dirtied</code>	<code>int8</code>	Number of database blocks dirtied by vacuum operations performed on this database
<code>total_blks_written</code>	<code>int8</code>	Number of database blocks written by vacuum operations performed on this database
<code>wal_records</code>	<code>int8</code>	Total number of WAL records generated by vacuum operations performed on this database
<code>wal_fpi</code>	<code>int8</code>	Total number of WAL full page images generated by vacuum operations performed on this database
<code>wal_bytes</code>	<code>numeric</code>	Total amount of WAL bytes generated by vacuum operations performed on this database
<code>blk_read_time</code>	<code>float8</code>	Time spent reading database blocks by vacuum operations performed on this database, in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>blk_write_time</code>	<code>float8</code>	Time spent writing database blocks by vacuum operations performed on this database, in milliseconds (if track_io_timing is enabled, otherwise zero)
<code>delay_time</code>	<code>float8</code>	Time spent sleeping in a vacuum delay point by vacuum operations performed on this database, in milliseconds (see Section 19.4.4 for details)
<code>system_time</code>	<code>float8</code>	System CPU time of vacuuming this database, in milliseconds
<code>user_time</code>	<code>float8</code>	User CPU time of vacuuming this database, in milliseconds
<code>total_time</code>	<code>float8</code>	Total time of vacuuming this database, in milliseconds

Column	Type	Description
interrupts	int4	Number of times vacuum operations performed on this database were interrupted on any errors

F.7.4.8. The `pgpro_stats_vacuum_tables` View

Important

Starting with Postgres Pro 16, this view contains no data because the statistics to be displayed are available through the catalog view `pg_stats_vacuum_tables` (see [System Views](#) for details).

The `pgpro_stats_vacuum_tables` view will contain one row for each table in the current database (including TOAST tables), showing statistics about vacuuming that specific table. These statistics are collected by the core system as explained in [Section 27.2](#). The table below describes the columns of the view.

Table F.15. `pgpro_stats_vacuum_tables` Columns

Column	Type	Description
relid	oid	OID of a table
schema	name	Name of the schema this table is in
relname	name	Name of this table
total_blks_read	int8	Number of database blocks read by vacuum operations performed on this table
total_blks_hit	int8	Number of times database blocks were found in the buffer cache by vacuum operations performed on this table
total_blks_dirtied	int8	Number of database blocks dirtied by vacuum operations performed on this table
total_blks_written	int8	Number of database blocks written by vacuum operations performed on this table
rel_blks_read	int8	Number of blocks vacuum operations read from this table
rel_blks_hit	int8	Number of times blocks of this table were already found in the buffer cache by vacuum operations, so that a read was not necessary (this only includes hits in the Postgres Pro buffer cache, not the operating system's file system cache)
pages_scanned	int8	Number of pages examined by vacuum operations performed on this table
pages_removed	int8	Number of pages removed from the physical storage by vacuum operations performed on this table

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Column	Type	Description
pages_frozen	int8	Number of times vacuum operations marked pages of this table as all-frozen in the visibility map
pages_all_visible	int8	Number of times vacuum operations marked pages of this table as all-visible in the visibility map
tuples_deleted	int8	Number of dead tuples vacuum operations deleted from this table
tuples_frozen	int8	Number of tuples of this table that vacuum operations marked as frozen
dead_tuples	int8	Number of dead tuples vacuum operations left in this table due to their visibility in transactions
index_vacuum_count	int8	Number of times indexes on this table were vacuumed
rev_all_frozen_pages	int8	Number of times the all-frozen mark in the visibility map was removed for pages of this table
rev_all_visible_pages	int8	Number of times the all-visible mark in the visibility map was removed for pages of this table
wal_records	int8	Total number of WAL records generated by vacuum operations performed on this table
wal_fpi	int8	Total number of WAL full page images generated by vacuum operations performed on this table
wal_bytes	numeric	Total amount of WAL bytes generated by vacuum operations performed on this table
blk_read_time	float8	Time spent reading database blocks by vacuum operations performed on this table, in milliseconds (if track_io_timing is enabled, otherwise zero)
blk_write_time	float8	Time spent writing database blocks by vacuum operations performed on this table, in milliseconds (if track_io_timing is enabled, otherwise zero)
delay_time	float8	Time spent sleeping in a vacuum delay point by vacuum operations performed on this table, in milliseconds (see Section 19.4.4 for details)
system_time	float8	System CPU time of vacuuming this table, in milliseconds

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Column	Type	Description
user_time	float8	User CPU time of vacuuming this table, in milliseconds
total_time	float8	Total time of vacuuming this table, in milliseconds
interrupts	integer	Number of times vacuum operations performed on this table were interrupted on any errors

Columns `total_*`, `wal_*` and `blk_*` include data on vacuuming indexes on this table, while columns `system_time` and `user_time` only include data on vacuuming the heap.

F.7.4.9. The `pgpro_stats_vacuum_indexes` View

Important

Starting with Postgres Pro 16, this view contains no data because the statistics to be displayed are available through the catalog view `pg_stats_vacuum_indexes` (see [System Views](#) for details).

The `pgpro_stats_vacuum_indexes` view will contain one row for each index in the current database (including TOAST table indexes), showing statistics about vacuuming that specific index. These statistics are collected by the core system as explained in [Section 27.2](#). The table below describes the columns of the view.

Table F.16. `pgpro_stats_vacuum_indexes` Columns

Column	Type	Description
relid	oid	OID of an index
schema	name	Name of the schema this index is in
relname	name	Name of this index
total_blks_read	int8	Number of database blocks read by vacuum operations performed on this index
total_blks_hit	int8	Number of times database blocks were found in the buffer cache by vacuum operations performed on this index
total_blks_dirtied	int8	Number of database blocks dirtied by vacuum operations performed on this index
total_blks_written	int8	Number of database blocks written by vacuum operations performed on this index
rel_blks_read	int8	Number of blocks vacuum operations read from this index
rel_blks_hit	int8	Number of times blocks of this index were already found in the buffer cache by vacuum operations, so that a read was not necessary (this only includes hits in the Postgres Pro buffer cache, not the operating system's file system cache)

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Column	Type	Description
pages_deleted	int8	Number of pages deleted by vacuum operations performed on this index
tuples_deleted	int8	Number of dead tuples vacuum operations deleted from this index
wal_records	int8	Total number of WAL records generated by vacuum operations performed on this index
wal_fpi	int8	Total number of WAL full page images generated by vacuum operations performed on this index
wal_bytes	numeric	Total amount of WAL bytes generated by vacuum operations performed on this index
blk_read_time	float8	Time spent reading database blocks by vacuum operations performed on this index, in milliseconds (if track_io_timing is enabled, otherwise zero)
blk_write_time	float8	Time spent writing database blocks by vacuum operations performed on this index, in milliseconds (if track_io_timing is enabled, otherwise zero)
delay_time	float8	Time spent sleeping in a vacuum delay point by vacuum operations performed on this index, in milliseconds (see Section 19.4.4 for details)
system_time	float8	System CPU time of vacuuming this index, in milliseconds
user_time	float8	User CPU time of vacuuming this index, in milliseconds
total_time	float8	Total time of vacuuming this index, in milliseconds
interrupts	integer	Number of times vacuum operations performed on this index were interrupted on any errors

F.7.4.10. The pgpro_stats_sdm_stats_updated View

This view displays the number of statistics entries received from each shard node and the timestamp of the last received statistics. It contains one row for each replication group from which statistics on query fragments were received. The columns of the view are shown in [Table F.17](#).

Table F.17. pgpro_stats_sdm_stats_updated Columns

Name	Type	Description
rgid	integer	ID of the replication group from which statistics with query fragments were received

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
updated	timestamp	Time at which the statistics on query fragments were last received from this replication group
num_entries	integer	Total number of statistics entries on query fragments received from this replication group

F.7.4.11. Global Views

Postgres Pro Shardman has global views based on the Postgres Pro local views. The definition of global view columns is the same as in its corresponding local view. Fetching from a global view returns a union of rows from the corresponding local views. The rows are fetched from each of their cluster nodes. Another difference is that the global views have an added column `rgid`. The `rgid` value shows the replication group ID of the cluster node from which a row is fetched. The `pgpro_stats` extension must be created on all the cluster nodes for the global views to work.

Table F.18. Global and Local Views

Global View	Local View
<code>shardman.gv_stats_statements</code>	<code>pgpro_stats_statements</code>
<code>shardman.gv_stats_sdm_statements</code>	<code>pgpro_stats_sdm_statements</code>

F.7.5. Data Types

F.7.5.1. The `pgpro_stats_rusage` Type

`pgpro_stats_rusage` is a record type that contains resource usage statistics of statement planning/execution. The fields of this type are shown in [Table F.19](#).

Table F.19. `pgpro_stats_rusage` Fields

Name	Type	Description
<code>reads</code>	<code>bigint</code>	Number of bytes read by the filesystem layer
<code>writes</code>	<code>bigint</code>	Number of bytes written by the filesystem layer
<code>user_time</code>	<code>double precision</code>	User CPU time used
<code>system_time</code>	<code>double precision</code>	System CPU time used
<code>minflts</code>	<code>bigint</code>	Number of page reclaims (soft page faults)
<code>majflts</code>	<code>bigint</code>	Number of page faults (hard page faults)
<code>nswaps</code>	<code>bigint</code>	Number of swaps
<code>msgsnds</code>	<code>bigint</code>	Number of IPC messages sent
<code>msgrcv</code>	<code>bigint</code>	Number of IPC messages received
<code>nsignals</code>	<code>bigint</code>	Number of signals received
<code>nvcsws</code>	<code>bigint</code>	Number of voluntary context switches
<code>nivcsws</code>	<code>bigint</code>	Number of involuntary context switches

F.7.6. Functions

`pgpro_stats_statements_reset(userid Oid, dbid Oid, queryid bigint, planid bigint, minmax_only boolean)` returns timestamp with time zone

`pgpro_stats_statements_reset` discards statistics gathered so far by `pgpro_stats` corresponding to the specified `userid`, `dbid`, `queryid`, and `planid`. If any of the parameters are not specified, the default value 0(invalid) is used for each of them and the statistics that match with other parameters will be reset. If no parameter is specified or all the specified parameters are 0(invalid), it will discard all statistics. If all statistics in the `pgpro_stats_statements` view are discarded, it will also reset the statistics in the `pgpro_stats_info` view. When `minmax_only` is true only the values of minimum and maximum planning and execution time will be reset (i.e. `min_plan_time`, `max_plan_time`, `min_exec_time` and `max_exec_time` fields). The default value for `minmax_only` parameter is false. Time of last min/max reset performed is shown in `minmax_stats_since` field of the `pgpro_stats_statements` view. This function returns the time of a reset. This time is saved to `stats_reset` field of `pgpro_stats_info` view or to `minmax_stats_since` field of the `pgpro_stats_statements` view if the corresponding reset was actually performed. By default, this function can only be executed by superusers. Access may be granted to others using `GRANT`.

Note

As statistics in the `pgpro_stats_vacuum_database`, `pgpro_stats_vacuum_tables`, and `pgpro_stats_vacuum_indexes` views are collected by the core system, to reset them, call the `pg_stat_reset()` function (see [Section 27.2.28](#) for details).

`pgpro_stats_statements(showtext boolean)` returns setof record

The `pgpro_stats_statements` view is defined in terms of a function also named `pgpro_stats_statements`. Users can also call the `pgpro_stats_statements` function directly, and by specifying `showtext := false` make query text be omitted (that is, the `OUT` argument that corresponds to the view's `query` column will return nulls). This feature is intended to support external tools that might wish to avoid the overhead of repeatedly retrieving query texts of indeterminate length. Such tools can instead cache the first query text observed for each entry themselves, since that is all `pgpro_stats` itself does, and then retrieve query texts only as needed. Since the server stores query texts in a file, this approach may reduce physical I/O for repeated examination of the `pgpro_stats_statements` data.

`pgpro_stats_info()` returns record

`pgpro_stats_info` view is defined in terms of a function also named `pgpro_stats_info`. Users can also call the `pgpro_stats_info` function directly.

`pgpro_stats_totals_reset(type text, id bigint)` returns timestamp with timezone

`pgpro_stats_totals_reset` discards statistics gathered so far by `pgpro_stats` corresponding to the specified object `type` and `id`. If no parameter is specified or the `type` parameter is set to 0, all statistics will be discarded. If `type` is set to a valid object type, then if `id` is specified, then statistics will be discarded only for the specified object, else, statistics will be discarded for all objects of the specified type. Otherwise, no statistics will be discarded. This function returns the time of a reset. By default, this function can only be executed by superusers. Access may be granted to others using `GRANT`.

`pgpro_stats_totals()` returns setof record

The `pgpro_stats_totals` view is defined in terms of a function also named `pgpro_stats_totals`. Users can also call the `pgpro_stats_totals` function directly.

`pgpro_stats_metrics()` returns setof record

Defines the `pgpro_stats_metrics` view, which is described in detail in [Table F.12](#).

`pgpro_stats_get_archiver()` returns setof record

Defines the `pgpro_stats_archiver` view, which is described in detail in [Table F.13](#).

`pgpro_stats_wal_sender_crc_errors()` returns bigint

In Postgres Pro Enterprise, returns the number of errors detected by the WAL sender process when the `wal_sender_check_crc` parameter is on. In Postgres Pro Standard, returns zero. Note that this counter is reset to zero when the cluster restarts.

`pgpro_stats_vacuum_database(dboid oid)` returns setof record

Defines the row of the `pgpro_stats_vacuum_database` view, which is described in detail in [Table F.14](#), for the database specified by `dboid`.

`pgpro_stats_vacuum_tables(dboid oid, relid oid)` returns setof record

Defines the row of the `pgpro_stats_vacuum_tables` view, which is described in detail in [Table F.15](#), for the database specified by `dboid` and table specified by `relid`. If `relid = 0`, the statistics for each table in the specified database are returned.

`pgpro_stats_vacuum_indexes(dboid oid, relid oid)` returns setof record

Defines the row of the `pgpro_stats_vacuum_indexes` view, which is described in detail in [Table F.16](#), for the database specified by `dboid` and index specified by `relid`. If `relid = 0`, the statistics for each index in the specified database are returned.

`pgpro_stats_sdm_stats_updated()` returns setof record

The `pgpro_stats_sdm_stats_updated` view is defined in terms of a function also named `pgpro_stats_sdm_stats_updated`. Users can also call the `pgpro_stats_sdm_stats_updated` function directly.

`pgpro_stats_sdm_stats_updated_reset()` returns void

Resets the number of statistics entries received from each shard node and the timestamp of the last received statistics.

F.7.6.1. Functions for Creating Views that Emulate Other Extensions

`pgpro_stats` can create views similar to those available in [pg_stat_statements](#) and [pg_stat_kcache](#) extensions. Specifically, `pg_stat_statements`, `pg_stat_statements_info`, `pg_stat_kcache` and `pg_stat_kcache_detail` views can be created. Each view is only created in a Postgres Pro version if it is available in `pg_stat_statements/pg_stat_kcache` extension for the same version of Postgres Pro/PostgreSQL. For example, the `pg_stat_statements_info` view is only created in Postgres Pro versions starting with 14. The following functions enable creating these views:

`pgpro_stats_create_pg_stat_statements_compatible_views()` returns void

Creates `pg_stat_statements` and `pg_stat_statements_info` views.

`pgpro_stats_create_pg_stat_kcache_compatible_views()` returns void

Creates `pg_stat_kcache` and `pg_stat_kcache_detail` views.

By default, these functions can only be executed by superusers. Access may be granted to others using `GRANT`.

To create `pg_stat_statements*` views, drop the `pg_stat_statements` extension if it was previously installed and call the function:

```
select pgpro_stats_create_pg_stat_statements_compatible_views();
```

To create `pg_stat_kcache*` views, drop the `pg_stat_kcache` extension if it was previously installed and call the function:

```
select pgpro_stats_create_pg_stat_kcache_compatible_views();
```


Once the views are created, you can work with them as if the `pg_stat_statements/pg_stat_kcache` extension is installed.

If you need to remove the views created earlier, do it in a regular way:

```
drop view pg_stat_statements;  
drop view pg_stat_statements_info;  
drop view pg_stat_kcache;  
drop view pg_stat_kcache_detail;
```

F.7.7. Configuration Parameters

F.7.7.1. General Settings

`pgpro_stats.max` (integer)

`pgpro_stats` only collects statistics about most frequent queries. The less frequently a query appears during the server operation, the less probable its inclusion in the statistics. Rare queries are almost immediately evicted from the statistics by more frequent queries. The `pgpro_stats.max` parameter defines the maximum number of unique pairs (normalized query text, normalized plan text) tracked by the module, (i.e., the maximum number of rows in the `pgpro_stats_statements` view). The larger this value, the larger the number of queries for which the information is stored. But this is achieved at the cost of reduced server performance when waiting for locks trying to access the statistics table in the shared memory and when periodically collecting garbage in the file with texts of queries and plans. The default value is 5000. This parameter can only be set at server start.

`pgpro_stats.max_totals` (integer)

`pgpro_stats.max_totals` is the maximum number of objects tracked by the module (i.e., the maximum number of rows in the `pgpro_stats_totals` view). If more distinct objects than that are observed, information about least-used objects is discarded. The default value is 1000. This parameter can only be set at server start.

`pgpro_stats.track` (enum)

`pgpro_stats.track` controls which statements are counted by the module. Specify `top` to track top-level statements (those issued directly by clients), `all` to also track nested statements (such as statements invoked within functions) with nesting level not greater than 100, or `none` to disable statement statistics collection. The default value is `top`. Only superusers can change this setting.

`pgpro_stats.track_utility` (boolean)

`pgpro_stats.track_utility` controls whether utility commands are tracked by the module. Utility commands are all those other than `SELECT`, `INSERT`, `UPDATE` and `DELETE`. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_planning` (boolean)

`pgpro_stats.track_planning` controls whether planning operations and duration are tracked by the module. Enabling this parameter may incur a noticeable performance penalty, especially when statements with identical query structure are executed by many concurrent connections which compete to update a small number of `pgpro_stats_statements` entries. The default value is `off`. Only superusers can change this setting.

`pgpro_stats.track_totals` (boolean)

`pgpro_stats.track_totals` controls whether aggregate statistics for objects (cluster, users, databases etc.) are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_cluster` (boolean)

`pgpro_stats.track_cluster` controls whether aggregate statistics for the cluster are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_databases` (boolean)

`pgpro_stats.track_databases` controls whether aggregate statistics for the databases are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_users` (boolean)

`pgpro_stats.track_users` controls whether aggregate statistics for the users are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_applications` (boolean)

`pgpro_stats.track_applications` controls whether aggregate statistics for the applications (whose names are set by [application_name](#)) are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_client_addr` (boolean)

`pgpro_stats.track_client_addr` controls whether aggregate statistics for the client IP addresses are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_backends` (boolean)

`pgpro_stats.track_backends` controls whether aggregate statistics for the backends are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.track_sessions` (boolean)

`pgpro_stats.track_sessions` controls whether aggregate statistics for the sessions are tracked by the module. The default value is `on`. Only superusers can change this setting.

`pgpro_stats.save` (boolean)

`pgpro_stats.save` specifies whether to save statement statistics across server shutdowns. If it is `off` then statistics are neither saved at shutdown nor reloaded at server start. The default value is `on`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

`pgpro_stats.plan_format` (text)

`pgpro_stats.plan_format` selects the `EXPLAIN` format for the query plan. Possible values are `text`, `xml`, `json`, and `yaml`. The default value is `text`. Changing this parameter requires a server restart.

`pgpro_stats.enable_profile` (boolean)

`pgpro_stats.enable_profile` enables sampling of wait events for separate statements. The default value is `true`. Changing this parameter requires a server restart.

`pgpro_stats.query_sample_rate` (float)

As `pgpro_stats` only collects statistics about most frequent queries, and rare queries are not included in the statistics, you can reduce the effect of the statistics collection on the server performance. This is possible because some time later than the statistics collection starts, almost all the statistics collected by the module will contain queries that are not very infrequent, and the set of types of the collected queries will not further change much if the value of [pgpro_stats.max](#) is sufficient and the average frequency for most of the frequent queries has no significant changes with the periods larger than the observation time. Such a stationary distribution of collected queries can also be achieved by increasing the time of the statistics collection, for example, twice, but only adding each second query to the statistics instead of every query, or considering each query, but adding queries to the statistics randomly, with a given probability. In this example, it equals 0.5. This allows you to considerably reduce the server load at the cost of increasing the time of the statistics collection. `pgpro_stats.query_sample_rate` is the parameter to implement this approach. For example: if `pgpro_stats.query_sample_rate` is set to 0.2 (1/5), in five hours of the server operation under

the above conditions, the statistics will be collected that is similar to the statistics collected with `pgpro_stats.query_sample_rate = 1` in an hour.

You can set up `pgpro_stats.query_sample_rate` by experiment. Choose the observation time depending on the environment of the server operation, for example: one hour. Look at the collected statistics an hour after the start and then an hour later without resetting the statistics. If the main difference appears to be in the increased number of collected queries (`calls`) rather than in the type of the queries, you can safely reduce `pgpro_stats.query_sample_rate` twice. In the case of a considerable difference in the types of the collected queries, increase the observation time and/or the `pgpro_stats.max` parameter.

The default value is 1.0. Changing this parameter requires a server restart.

`pgpro_stats.profile_period (integer)`

`pgpro_stats.profile_period` specifies the period, in milliseconds, during which to sample wait events. The default value is 10. Only superusers can change this setting.

`pgpro_stats.metrics_buffer_size (integer)`

`pgpro_stats.metrics_buffer_size` specifies the size of the ring buffer used for collecting statistical metrics. The default value is 16kB. Changing this parameter requires a server restart.

`pgpro_stats.metrics_workers (integer)`

`pgpro_stats.metrics_workers` specifies the number of workers used to collect statistical metrics. If this parameter is set to 2 or higher, one of the workers serves as the primary worker distributing queries to other workers. If only one worker is available, it gets reloaded to connect to different databases. Setting this parameter to 0 disables metrics collection. The default value is 2. Changing this parameter requires a server restart.

`pgpro_stats.stats_temp_directory (string)`

`pgpro_stats.stats_temp_directory` specifies the directory with the external file to store query texts. This can be a path relative to the data directory or an absolute path. Changing this parameter requires a server restart.

`pgpro_stats.track_sharded (boolean)`

`pgpro_stats.track_sharded` specifies whether the sharded statements are tracked and aggregated by `pgpro_stats`. The default value is `on`.

`pgpro_stats.pgss_max_nodes_tracked (integer)`

Sets the maximum number of nodes that are tracked by `pgpro_stats` for query fragments. It actually sets the maximum amount of the status entries that `pgpro_stats` can store for the `pgpro_stats_sdm_stats_updated` function. It does not affect the statistics tracking itself. The default value is 2048.

`pgpro_stats.transport_compression (string)`

Sets the algorithm for transport compression during statistics transferring between nodes. Transport compression is used to compress statistical entries passed from the shard nodes to the coordinator. Possible values are `pglz`, `zlib`, `lz4`, `zstd` or `off`. The default is `pglz`.

`pgpro_stats.enable_wait_counters (boolean)`

Enables or disables statistics collection for wait counters by enabling or disabling functions that calculate metrics of wait events. The default value is `off`.

`pgpro_stats.enable_rusage_counters (boolean)`

Enables or disables statistics collection for resource usage counters by enabling or disabling functions that calculate metrics of OS resource usage. The default value is `off`.

`pgpro_stats.track_shardman_connections` (enum)

Enables or disables statement processing for a distributed system. Possible values are `none` with no processing, `normalized` (default) with generalized statements being processed, and `all` with all statements being processed.

F.7.7.2. Metrics Settings

The following parameters can be used to define a custom metric to collect. The *N* placeholder in the parameter name serves as a unique identifier of the metric to which this setting should apply; it must be set to a non-negative integer for each metric.

When you add these parameters for a new metric, you have to restart the server for the changes to take effect. Once the new metric is added, its parameters can be changed without a server restart by simply reloading the `postgresql.conf` configuration file.

`pgpro_stats.metric_N_name` (text)

The name of metric *N*. This name will be displayed in the `metric_name` column of the `pgpro_stats_metrics` view.

`pgpro_stats.metric_N_query` (text)

The query statement that defines the metric to collect.

`pgpro_stats.metric_N_period` (integer)

The time interval at which to collect metric *N*, in milliseconds. Default: 60000 ms

`pgpro_stats.metric_N_db` (text)

The list of databases for which to collect metric *N*. Database names must be separated by commas. You can specify the `*` character to select all databases in the cluster except the template databases. If you need to analyze the template databases as well, you have to specify them explicitly.

`pgpro_stats.metric_N_user` (text)

The name of the user on behalf of which to collect metric *N*. This user must have access to the database for which the metric is collected.

F.7.8. Tracing of Application Sessions

In `pgpro_stats`, tracing of application sessions is implemented. It is based on filters, which trigger logging the execution of queries that match filtering conditions. Queries and their `EXPLAIN` output are logged in so-called *trace files* specified by the user or in the system log file (if the trace file is not specified). Filters are stored in a table located in the shared memory. The rows of this table represent filters, and the columns contain filtering conditions and `EXPLAIN` options to be used (see [EXPLAIN](#) for details). You should fill this table with filters to start tracing of queries.

Once a database administrator adds a filter in any session, all subsequent executions of queries that match the filter conditions will be traced by all sessions of the instance without a need in the server restart. In other words, filters can be added, deleted, or updated “on the fly”, and tracing with these filters immediately starts for existing and future sessions.

Note

Operations of adding, deleting, or updating filters are not transactional. Changes to the table of filters are not rolled back together with the transaction where these changes were made.

Columns of the table that stores filters are also called *filter attributes*. They are described in [Table F.20](#).

Table F.20. Filter Attributes

Name	Type	Description
filter_id	integer	Unique filter ID, which is automatically assigned to a filter when it is created, numbered from 1.
active	boolean	With the value of <code>false</code> , the filter is not active, which means that queries matching this filter are not logged. Default: <code>true</code> .
alias	name	Filter name, which can be specified for convenience. Only shown in the output of the <code>pgpro_stats_trace_show()</code> function. Default: empty string.
tracefile	name	Name of the trace file, where queries matching the filter and their <code>EXPLAIN</code> output are logged. If this filename is not specified, logging is done to the system log with the log level specified by <code>pgpro_stats.trace_log_level</code> . Trace files are created in <code>PGDATA/pg_stat</code> directory and have <code>trace</code> extension.
pid	integer	Process ID of the backend that executes the statement
database_name	name	Name of the database where the statement is executed
client_addr	name	IP address of the client connected to this backend
application_name	name	Name of the application that invoked execution of the statement
username	name	Name of the user who executes the statement
queryid	bigint	Internal hash code, computed from the statement's parse tree
planid	bigint	Internal hash code, computed from the statement's plan tree
duration	float8	Time spent in the planning and execution of the statement, in milliseconds
plan_time	float8	Time spent in the planning of the statement, in milliseconds
exec_time	float8	Time spent in the execution of the statement, in milliseconds
user_time	float8	User CPU time used in planning and execution of the statement
system_time	float8	System CPU time used in planning and execution of the statement

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
rows	int8	Total number of rows retrieved or affected by the statement
shared_blks_hit	int8	Total number of shared block cache hits by the statement
shared_blks_read	int8	Total number of shared blocks read by the statement
shared_blks_fetched	int8	Total number of shared blocks fetched from buffers by the statement
shared_blks_dirtied	int8	Total number of shared blocks dirtied by the statement
shared_blks_written	int8	Total number of shared blocks written by the statement
local_blks_hit	int8	Total number of local block cache hits by the statement
local_blks_read	int8	Total number of local blocks read by the statement
local_blks_fetched	int8	Total number of local blocks fetched from buffers by the statement
local_blks_dirtied	int8	Total number of local blocks dirtied by the statement
local_blks_written	int8	Total number of local blocks written by the statement
temp_blks_read	int8	Total number of temp blocks read by the statement
temp_blks_written	int8	Total number of temp blocks written by the statement
wal_bytes	numeric	Total amount of WAL bytes generated by the statement
total_wait_time	float8	Total time execution of this statement spent waiting
total_inval_msgs	bigint	Total number of cache invalidation messages generated by the statement (if this is supported by the server)
explain_analyze	boolean	If true, the <code>EXPLAIN</code> output will be logged with the <code>ANALYZE</code> option. By default, defined by the pgpro_stats.explain_analyze_default configuration parameter.
explain_verbose	boolean	If true, the <code>EXPLAIN</code> output will be logged with the <code>VERBOSE</code> option. By default, defined by the pgpro_stats.explain_verbose_default configuration parameter.
explain_costs	boolean	If true, the <code>EXPLAIN</code> output will be logged with the <code>COSTS</code> option. By default, defined by the pgpro_

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
		stats.explain_costs_default configuration parameter.
explain_settings	boolean	If true, the EXPLAIN output will be logged with the SETTINGS option. By default, defined by the pgpro_stats.explain_settings_default configuration parameter.
explain_buffers	boolean	If true, the EXPLAIN output will be logged with the BUFFERS option. By default, defined by the pgpro_stats.explain_buffers_default configuration parameter.
explain_wal	boolean	If true, the EXPLAIN output will be logged with the WAL option. By default, defined by the pgpro_stats.explain_wal_default configuration parameter.
explain_timing	boolean	If true, the EXPLAIN output will be logged with the TIMING option. By default, defined by the pgpro_stats.explain_timing_default configuration parameter.
explain_format	text	The value of the FORMAT option of EXPLAIN to be logged, which can be TEXT, XML, JSON, or YAML. By default, defined by the pgpro_stats.plan_format configuration parameter.
time_info	boolean	If true, includes additional information in the session-tracing output in the format defined by <code>explain_format : pid, duration</code> in milliseconds, start time — start time of the query execution, and stop time — end time of the query execution. See Example F.6 for details.

Filter attributes can be divided into the following groups:

- Control attributes, from `filter_id` to `tracefile` in [Table F.20](#).
- Identification attributes, from `pid` to `planid` in [Table F.20](#).
- Threshold attributes, from `duration` to `total_inval_msgs` in [Table F.20](#).
- Attributes that only control the output to the trace file or system log file and do not affect filtering. Most of them specify the EXPLAIN options. These are the attributes starting with `explain_analyze` up to the end of [Table F.20](#).

Query execution will be traced if both of these conditions are met:

- For all the specified identification attributes of the filter, the values of the respective characteristics of the query are the same as the values of these attributes.
- For all the specified threshold attributes of the filter, the values of the respective characteristics of the query are greater than or equal to the limits specified in these attributes.

It is allowed not to assign values to attributes. For all the attributes except `active`, you can also explicitly assign the value of `NULL`. If the value of a control attribute is undefined or it was explicitly set to `NULL`, this attribute will be assigned the default value shown in [Table F.20](#). For identification and threshold attributes, undefined or `NULL` values mean that queries will not be filtered by the respective characteristic. If no values are defined for all identification and all threshold attributes, all the queries are taken to be matching this filter and will be logged. If the attributes specifying `EXPLAIN` options are undefined, the default values defined by the configuration parameters will be used (see [Section F.7.8.2](#) for details).

Warning

Although when specifying a filter, you can assign values to any combination of filter attributes, bear in mind that a too general filter will lead to an excessive size of the trace file and will affect the performance more than desired as the main performance overhead is associated with writing to the trace file rather than with checking the filter conditions.

F.7.8.1. Session-Tracing Functions

Specialized functions enable creation, update and deletion of query filters:

`pgpro_stats_trace_insert(VARIADIC "any")` returns integer

Adds a filter to the list of session-tracing filters. A filter must be passed as a sequence of alternating key-value pairs, where the key is the name of the filter attribute from [Table F.20](#) except `filter_id`. For example:

```
pgpro_stats_trace_insert('pid', 42, 'database_name', 'main', 'explain_analyze',
true)
```

Key-value pairs are position-independent: they can be provided in any order, independently of the order of attributes in [Table F.20](#). Some of the attributes can be missing from the sequence. Such attributes will be assigned the default values. Returns the unique `filter_id`, which the function assigns to the added filter.

`pgpro_stats_trace_update(filter_id integer, VARIADIC "any")` returns boolean

Updates a session-tracing filter defined by `filter_id`. Filter attributes to update must be passed as a sequence of alternating key-value pairs, where the key is the name of the filter attribute from [Table F.20](#) except `filter_id`. Key-value pairs are position-independent: they can be provided in any order, independently of the order of attributes in [Table F.20](#). Filter attributes whose values are not explicitly specified will not be changed. The value of `NULL` for identification and threshold attributes resets the values of these attributes, which means that a query will match the conditions of the updated filter regardless of the respective characteristics. Returns `true` on success and `false` otherwise.

`pgpro_stats_trace_delete(filter_id integer)` returns boolean

Deletes a session-tracing filter defined by `filter_id`. Returns `true` on success and `false` otherwise.

`pgpro_stats_trace_reset()` returns integer

Removes all session-tracing filters. Returns the number of removed filters. By default, this function can only be executed by superusers. Access may be granted to others using `GRANT`.

`pgpro_stats_trace_show()` returns setof record

Displays the contents of the table of filters, whose columns are shown in [Table F.20](#).

F.7.8.2. Session-Tracing Configuration Parameters

The following parameters can be used to configure session-tracing logging. Note that `pgpro_stats.explain_*_default` settings define the logging behavior when the corresponding `explain_*` filter attributes (from [Table F.20](#)) are not specified or explicitly set to `NULL`.

`pgpro_stats.explain_analyze_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `ANALYZE` option and without it if `false`. The default value is `false`. Only superusers can change this setting.

`pgpro_stats.explain_verbose_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `VERBOSE` option and without it if `false`. The default value is `false`. Only superusers can change this setting.

`pgpro_stats.explain_costs_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `COSTS` option and without it if `false`. The default value is `true`. Only superusers can change this setting.

`pgpro_stats.explain_settings_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `SETTINGS` option and without it if `false`. The default value is `false`. Only superusers can change this setting.

`pgpro_stats.explain_buffers_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `BUFFERS` option and without it if `false`. The default value is `false`. Only superusers can change this setting.

`pgpro_stats.explain_wal_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `WAL` option and without it if `false`. The default value is `false`. Only superusers can change this setting.

`pgpro_stats.explain_timing_default` (boolean)

If `true`, the `EXPLAIN` output will be logged with the `TIMING` option and without it if `false`. The default value is `true`. Only superusers can change this setting.

`pgpro_stats.trace_log_level` (enum)

Defines the log level at which the `EXPLAIN` output will be logged to the system log file (when the trace file is not specified). Valid values are `DEBUG5`, `DEBUG4`, `DEBUG3`, `DEBUG2`, `DEBUG1`, `INFO`, `NOTICE`, `WARNING`, and `LOG`. The default value is `LOG`.

`pgpro_stats.trace_query_text_size` (int)

Defines the maximum size of query texts in the session-tracing output, in kilobytes. If the size of the query text in kilobytes is larger, it is trimmed to the maximum size. If this parameter is not specified, the default value of zero is used, which means that the whole query texts are output.

F.7.8.3. Examples Related to Session Tracing

Example F.5. Usage of Session-Tracing Functions

Let's add the filter first:

```
SELECT pgpro_stats_trace_insert('alias', 'first', 'pid', pg_backend_pid(),  
    'explain_analyze', true);
```

Let's add the filter second and specify logging to the trace file `second_tf.trace`:

```
SELECT pgpro_stats_trace_insert('alias', 'second', 'database_name', current_database(),  
    'explain_costs', false, 'tracefile', 'second_tf');
```

You can view the table with filters as follows:

```
\x auto
```


Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

```
SELECT * from pgpro_stats_trace_show();
```

```
-[ RECORD 1 ]-----+-----
```

filter_id	1
active	t
alias	first
tracefile	
pid	243183
database_name	
client_addr	
application_name	
username	
queryid	
planid	
duration	
plan_time	
exec_time	
user_time	
system_time	
rows	
shared_blks_hit	
shared_blks_read	
shared_blks_fetched	
shared_blks_dirtied	
shared_blks_written	
local_blks_hit	
local_blks_read	
local_blks_fetched	
local_blks_dirtied	
local_blks_written	
temp_blks_read	
temp_blks_written	
wal_bytes	
total_wait_time	
total_inval_msgs	
explain_analyze	t
explain_verbose	f
explain_costs	t
explain_settings	f
explain_buffers	f
explain_wal	f
explain_timing	t
explain_format	text
time_info	f

```
-[ RECORD 2 ]-----+-----
```

filter_id	2
active	t
alias	second
tracefile	second_tf
pid	
database_name	postgres
client_addr	
application_name	
username	
queryid	
planid	
duration	
plan_time	

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

exec_time	
user_time	
system_time	
rows	
shared_blks_hit	
shared_blks_read	
shared_blks_fetched	
shared_blks_dirtied	
shared_blks_written	
local_blks_hit	
local_blks_read	
local_blks_fetched	
local_blks_dirtied	
local_blks_written	
temp_blks_read	
temp_blks_written	
wal_bytes	
total_wait_time	
total_inval_msgs	
explain_analyze	f
explain_verbose	f
explain_costs	f
explain_settings	f
explain_buffers	f
explain_wal	f
explain_timing	f
explain_format	text
time_info	f

The following query matches the conditions of both filters, so it must be logged in the system log file and in the specified trace file:

```
SELECT 1 as result;
```

The following is the output to the system log file:

```
2023-04-18 04:52:53.242 MSK [63112] LOG:  Filter 1 triggered explain of the plan:
Query Text: SELECT 1 as result;
Result  (cost=0.00..0.01 rows=1 width=4) (actual time=0.003..0.003 rows=1 loops=1)
```

And the following is the output to the `second_tf.trace` trace file:

```
Query Text: SELECT 1 as result;
Result
```

Let's delete the first filter:

```
SELECT pgpro_stats_trace_delete(1);
```

Let's also change `pid` to 2 for the second filter

```
SELECT pgpro_stats_trace_update(2, 'pid', 2);
```

When you execute the query

```
SELECT 2 as result;
```

it does not get logged to `second_tf.trace`.

Now let's remove all filters from the table:

```
SELECT pgpro_stats_trace_reset();
```

Example F.6. Session-Tracing Output with `time_info = true`

The following is an example of the session-tracing output in the JSON format when `time_info = true`:

```
{
```

```
"Pid": 123456,
"Duration ms": 0.094691,
"Start time": "2025-01-12 15:39:27.580447+03",
"Stop time": "2025-01-12 15:39:27.587802+03",
"Query Text": "SELECT 2 AS Result;",
"Plan": {
  "Node Type": "Result",
  "Parallel Aware": false,
  "Async Capable": false,
  "Startup Cost": 0.00,
  "Total Cost": 0.01,
  "Plan Rows": 1,
  "Plan Width": 4
}
}
```

F.7.9. Cache Invalidation Metrics

Among the rest, `pgpro_stats` can collect cache invalidation statistics. This section provides some background information needed to better understand related metrics.

Each backend has its local cache, which allows you to minimize accesses for meta information on tables, for example, to the system catalogs. If a backend changes the meta information, this information must be updated in other backends' caches. This is implemented by sending invalidation messages through a queue: the backend that changed the meta information on some object sends an appropriate message to the queue.

All backends get invalidation messages from the queue. Depending on whether the object for which the invalidation message was received is cached, the backend either ignores the message (when the object is not cached) or updates its cache (when the object is cached). In `pgpro_stats`, most invalidation message counters, unless explicitly stated otherwise for certain counters, are incremented when backends just generate messages, which will only be sent to the queue upon commit of the appropriate transaction. Note that the counters will remain incremented if the transaction is rolled back, although the message will not be sent to the queue.

When a backend that is adding messages to the queue figures out that the queue size reached a certain limit, it starts a cleanup by deleting messages already processed by all backends, and if backends are found that heavily fall behind and thus delay the cleanup, they get a reset signal, which forces them to reset all their caches.

F.7.9.1. The `pgpro_stats_inval_status` View

The `pgpro_stats_inval_status` view shows one row with the current status of the cache invalidation global queue. The columns of the view are shown in [Table F.21](#).

Table F.21. `pgpro_stats_inval_status` Columns

Name	Type	Description
<code>num_inval_messages</code>	<code>int8</code>	Current number of invalidation messages in the queue
<code>num_inval_queue_cleanups</code>	<code>int8</code>	Number of invalidation queue cleanups done to prevent its overflow
<code>num_inval_queue_resets</code>	<code>int4</code>	Number of cache resets for backends that fail to process messages fast enough

In a working system, `num_inval_messages` usually approximately equals 4000, which means that the queue is pretty full. The speed of the `num_inval_queue_cleanups` growth is determined by how fast in-

validation messages are generated. Growth of `num_inval_queue_resets` is normally zero, and non-zero growth indicates either too fast generation of messages or delays in processing messages by backends. Monitoring `num_inval_queue_cleanups` and `num_inval_queue_resets` may in some cases allow you to detect problematic backend/backends as described below.

If for a certain time interval, `num_inval_queue_cleanups` considerably increased, while `num_inval_queue_resets` did not, this indicates that invalidation messages are generated faster and/or backends process them more slowly, but backends still manage to process messages before the queue overflows.

If for a time interval, `num_inval_queue_cleanups` did not considerably increase, while `num_inval_queue_resets` did, this definitely indicates a delay in processing messages by backend(s), and the `cache_resets` column of the `pgpro_stats_totals` view allows you to figure out which backend(s) to blame.

If for a time interval, both counters considerably increased, this also indicates that invalidation messages are generated faster and/or backends process them more slowly, but this time backends fail to process messages before the queue overflows. The `cache_resets` column of the `pgpro_stats_totals` view allows you detect which backend(s) delay message processing. In this case, it is not possible to definitely conclude whether too fast generation of messages or a delay in message processing accounts for the growth of `num_inval_queue_resets`. However, the `totals` counter of the `pgpro_stats_inval_msgs` view may help here. If the change of this counter for that interval is pretty the same as for a previous interval of the same length, you can definitely conclude that the growth is caused by backend delays.

The `pgpro_stats_inval_status` view can be defined in terms of a function:

```
pgpro_stats_inval_status() returns record
```

Defines the `pgpro_stats_inval_status` view, which is described in detail in [Table F.21](#).

F.7.9.2. The `pgpro_stats_inval_msgs` Type

The `pgpro_stats_statements` and `pgpro_stats_totals` views for each corresponding object, show a record of the `pgpro_stats_inval_msgs` record type containing counters for cache invalidation messages. The fields of the type are shown in [Table F.22](#).

Table F.22. `pgpro_stats_inval_msgs` Fields

Name	Type	Description
<code>total</code>	<code>bigint</code>	Total number of invalidation messages
<code>catcache</code>	<code>bigint</code>	Number of selective catalog cache invalidation messages
<code>catalog</code>	<code>bigint</code>	Number of whole catalog cache invalidation messages
<code>relcache</code>	<code>bigint</code>	Number of selective relation cache invalidation messages
<code>relcache_all</code>	<code>bigint</code>	Number of whole relation cache invalidation messages
<code>smgr</code>	<code>bigint</code>	Number of invalidation messages of open relation files. Gets incremented when the messages are sent to the queue.
<code>relmap</code>	<code>bigint</code>	Number of relation map cache invalidation messages. Gets incremented when the messages are sent to the queue.

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Name	Type	Description
snapshot	bigint	Number of catalog snapshot invalidation messages

F.7.9.3. Related Configuration Parameters

`pgpro_stats.enable_inval_msgs_counters` (boolean)

Enables or disables statistics collection of the invalidation messages by enabling or disabling functions that calculate metrics of invalidation messages. If disabled, the `pgpro_stats_inval_status` view is empty. The default value is `off`.

F.7.10. Authors

Postgres Professional, Moscow, Russia

F.8. pg_query_state — a facility to know the current state of query execution on working backend

The `pg_query_state` module provides facility to know the current state of query execution on working backend and `silkworm` multiplexer workers.

F.8.1. Overview

Each non-utility query statement (SELECT/INSERT/UPDATE/DELETE) after optimization/planning stage is translated into plan tree, which is kind of imperative representation of declarative SQL query. EXPLAIN ANALYZE request allows to demonstrate execution statistics gathered from each node of plan tree (full time of execution, number of rows emitted to upper nodes, etc). But this statistics is collected after execution of query. This module allows to show actual statistics of query running on external backend. At that, format of resulting output is almost identical to ordinal EXPLAIN ANALYZE. Thus users are able to track of query execution in progress. In fact, this module is able to explore external backend and determine its actual state. Particularly it's helpful when backend executes a heavy query or gets stuck.

F.8.2. Use cases

Using this module there can help in the following things:

- detect a long query (along with other monitoring tools);
- supervise query execution.

F.8.3. Installation

To install `pg_query_state` run in `psql`:

```
CREATE EXTENSION pg_query_state;
```

Then modify `shared_preload_libraries` parameter in `postgres.conf` as following:

```
shared_preload_libraries = 'pg_query_state'
```

It will require to restart the Postgres Pro instance.

F.8.4. Function `pg_query_state`

```
pg_query_state(integer pid,  
               verbose boolean DEFAULT FALSE,  
               costs   boolean DEFAULT FALSE,  
               timing   boolean DEFAULT FALSE,  
               buffers  boolean DEFAULT FALSE,  
               triggers boolean DEFAULT FALSE,  
               format   text    DEFAULT 'text')  
returns TABLE (pid integer,  
                frame_number integer,  
                query_text text,  
                plan text,  
                leader_pid integer)
```

Extract current query state from backend with specified `pid`. Since parallel query can spawn multiple workers and function call causes nested subqueries so that state of execution may be viewed as stack of running queries, return value of `pg_query_state` has type `TABLE (pid integer, frame_number integer, query_text text, plan text, leader_pid integer)`. It represents tree structure consisting of leader

process and its spawned workers identified by `pid`. Each worker refers to leader through `leader_pid` column. For leader process the value of this column is `null`. The state of each process is represented as stack of function calls. Each frame of that stack is specified as correspondence between `frame_number` starting from zero, `query_text` and `plan` with online statistics columns.

Thus, user can see the states of main query and queries generated from function calls for leader process and all workers spawned from it.

In process of execution some nodes of plan tree can take loops of full execution. Therefore statistics for each node consists of two parts: average statistics for previous loops just like in `EXPLAIN ANALYZE` output and statistics for current loop if node have not finished.

Optional arguments:

- `verbose` - use `EXPLAIN VERBOSE` for plan printing;
- `costs` - costs for each node;
- `timing` - print timing data for each node, if collecting of timing statistics is turned off on called side resulting output will contain WARNING message `timing statistics disabled`;
- `buffers` - print buffers usage, if collecting of buffers statistics is turned off on called side resulting output will contain WARNING message `buffers statistics disabled`;
- `triggers` - include triggers statistics in result plan trees;
- `format` - `EXPLAIN` format to be used for plans printing, possible values: `text`, `xml`, `json`, `yaml`.

If callable backend is not executing any query the function prints INFO message about backend's state taken from `pg_stat_activity` view if it exists there.

Calling role have to be superuser or member of the role whose backend is being called. Otherwise function prints ERROR message `permission denied`.

F.8.5. Configuration settings

There are several user-accessible GUC variables designed to toggle the whole module and the collecting of specific statistic parameters while query is running:

- `pg_query_state.enable` - disable (or enable) `pg_query_state` completely, default value is `true`
- `pg_query_state.enable_timing` - collect timing data for each node, default value is `false`
- `pg_query_state.enable_buffers` - collect buffers usage, default value is `false`

These parameters are set on called side before running any queries whose states are attempted to extract. WARNING: if `pg_query_state.enable_timing` is turned off the calling side cannot get time statistics, similarly for `pg_query_state.enable_buffers` parameter.

F.8.6. Examples

Set maximum number of parallel workers on `gather` node equals 2:

```
postgres=# set max_parallel_workers_per_gather = 2;
```

Assume one backend with `pid = 49265` performs a simple query:

```
postgres=# select pg_backend_pid();
pg_backend_pid
-----
 49265
(1 row)
postgres=# select count(*) from foo join bar on foo.c1=bar.c1;
```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Other backend can extract intermediate state of execution that query:

```
postgres=# \x
postgres=# select * from pg_query_state(49265);
-[ RECORD
  1 ]+-----
pid          | 49265
frame_number | 0
query_text   | select count(*) from foo join bar on foo.c1=bar.c1;
plan         | Finalize Aggregate (Current loop: actual rows=0, loop number=1)
              +
              | -> Gather (Current loop: actual rows=0, loop number=1)
              +
              |       Workers Planned: 2
              +
              |       Workers Launched: 2
              +
              |       -> Partial Aggregate (Current loop: actual rows=0, loop
number=1)
              +
              |       -> Nested Loop (Current loop: actual rows=12, loop
number=1)
              +
              |               Join Filter: (foo.c1 = bar.c1)
              +
              |               Rows Removed by Join Filter: 5673232
              +
              |               -> Parallel Seq Scan on foo (Current loop: actual
rows=12, loop number=1)
              +
              |               -> Seq Scan on bar (actual rows=500000 loops=11)
              (Current loop: actual rows=173244, loop number=12)
leader_pid   | (null)
-[ RECORD
  2 ]+-----
pid          | 49324
frame_number | 0
query_text   | <parallel query>
plan         | Partial Aggregate (Current loop: actual rows=0, loop number=1)
              +
              | -> Nested Loop (Current loop: actual rows=10, loop number=1)
              +
              |       Join Filter: (foo.c1 = bar.c1)
              +
              |       Rows Removed by Join Filter: 4896779
              +
              |       -> Parallel Seq Scan on foo (Current loop: actual rows=10, loop
number=1)
              +
              |       -> Seq Scan on bar (actual rows=500000 loops=9) (Current loop:
actual rows=396789, loop number=10)
leader_pid   | 49265
-[ RECORD
  3 ]+-----
pid          | 49323
frame_number | 0
query_text   | <parallel query>
plan         | Partial Aggregate (Current loop: actual rows=0, loop number=1)
              +
              | -> Nested Loop (Current loop: actual rows=11, loop number=1)
              +
```

**Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib**

```

| Join Filter: (foo.c1 = bar.c1)
| +
| Rows Removed by Join Filter: 5268783
| +
| -> Parallel Seq Scan on foo (Current loop: actual rows=11, loop
number=1)
| +
| -> Seq Scan on bar (actual rows=500000 loops=10) (Current loop:
actual rows=268794, loop number=11)
leader_pid | 49265

```

In example above working backend spawns two parallel workers with pids 49324 and 49323. Their leader_pid column's values clarify that these workers belong to the main backend. Seq Scan node has statistics on passed loops (average number of rows delivered to Nested Loop and number of passed loops are shown) and statistics on current loop. Other nodes has statistics only for current loop as this loop is first (loop number = 1).

Assume first backend executes some function:

```
postgres=# select n_join_foo_bar();
```

Other backend can get the follow output:

```

postgres=# select * from pg_query_state(49265);
-[ RECORD
  1 ]+-----
pid          | 49265
frame_number | 0
query_text   | select n_join_foo_bar();
plan         | Result (Current loop: actual rows=0, loop number=1)
leader_pid   | (null)
-[ RECORD
  2 ]+-----
pid          | 49265
frame_number | 1
query_text   | SELECT (select count(*) from foo join bar on foo.c1=bar.c1)
plan         | Result (Current loop: actual rows=0, loop number=1)
| +
| InitPlan 1 (returns $0)
| +
| -> Aggregate (Current loop: actual rows=0, loop number=1)
| +
| -> Nested Loop (Current loop: actual rows=51, loop number=1)
| +
| Join Filter: (foo.c1 = bar.c1)
| +
| Rows Removed by Join Filter: 51636304
| +
| -> Seq Scan on bar (Current loop: actual rows=52, loop
number=1)
| +
| -> Materialize (actual rows=1000000 loops=51) (Current
loop: actual rows=636355, loop number=52)+
| -> Seq Scan on foo (Current loop: actual
rows=1000000, loop number=1)
leader_pid   | (null)

```

First row corresponds to function call, second - to query which is in the body of that function.

We can get result plans in different format (e.g. json):

Additional Supplied Modules and Extensions Shipped in

postgrespro-sdm-17-contrib

```

postgres=# select * from pg_query_state(pid := 49265, format := 'json');
-[ RECORD 1 ]+-----
pid          | 49265
frame_number | 0
query_text   | select * from n_join_foo_bar();
plan         | {
              |   "Plan": {
              |     "Node Type": "Function Scan",
              |     "Parallel Aware": false,
              |     "Function Name": "n_join_foo_bar",
              |     "Alias": "n_join_foo_bar",
              |     "Current loop": {
              |       "Actual Loop Number": 1,
              |       "Actual Rows": 0
              |     }
              |   }
              | }
leader_pid   | (null)
-[ RECORD 2 ]+-----
pid          | 49265
frame_number | 1
query_text   | SELECT (select count(*) from foo join bar on foo.c1=bar.c1)
plan         | {
              |   "Plan": {
              |     "Node Type": "Result",
              |     "Parallel Aware": false,
              |     "Current loop": {
              |       "Actual Loop Number": 1,
              |       "Actual Rows": 0
              |     },
              |     "Plans": [
              |       {
              |         "Node Type": "Aggregate",
              |         "Strategy": "Plain",
              |         "Partial Mode": "Simple",
              |         "Parent Relationship": "InitPlan",
              |         "Subplan Name": "InitPlan 1 (returns $0)",
              |         "Parallel Aware": false,
              |         "Current loop": {
              |           "Actual Loop Number": 1,
              |           "Actual Rows": 0
              |         },
              |       },
              |       "Plans": [
              |         {
              |           "Node Type": "Nested Loop",
              |           "Parent Relationship": "Outer",
              |           "Parallel Aware": false,
              |           "Join Type": "Inner",
              |           "Current loop": {
              |             "Actual Loop Number": 1,
              |             "Actual Rows": 610
              |           },
              |         },
              |         "Join Filter": "(foo.c1 = bar.c1)",
              |         "Rows Removed by Join Filter": 610072944,
              |         "Plans": [
              |           {
              |             "Node Type": "Seq Scan",
              |             "Parent Relationship": "Outer",

```

postgrespro-sdm-17-contrib

```
leader_pid
```

F.9. pg_stat_statements — track statistics of SQL planning and execution

The `pg_stat_statements` module provides a means for tracking planning and execution statistics of all SQL statements executed by a server.

The module must be loaded by adding `pg_stat_statements` to [shared_preload_libraries](#) in `postgresql.conf`, because it requires additional shared memory. This means that a server restart is needed to add or remove the module. In addition, query identifier calculation must be enabled in order for the module to be active, which is done automatically if [compute_query_id](#) is set to `auto` or `on`, or any third-party module that calculates query identifiers is loaded.

When `pg_stat_statements` is active, it tracks statistics across all databases of the server. To access and manipulate these statistics, the module provides views `pg_stat_statements` and `pg_stat_statements_info`, and the utility functions `pg_stat_statements_reset` and `pg_stat_statements`. These are not available globally but can be enabled for a specific database with `CREATE EXTENSION pg_stat_statements`.

F.9.1. The pg_stat_statements View

The statistics gathered by the module are made available via a view named `pg_stat_statements`. This view contains one row for each distinct combination of database ID, user ID, query ID and whether it's a top-level statement or not (up to the maximum number of distinct statements that the module can track). The columns of the view are shown in [Table F.23](#).

Table F.23. pg_stat_statements Columns

Column Type	Description
<code>userid oid</code> (references <code>pg_authid .oid</code>)	OID of user who executed the statement
<code>dbid oid</code> (references <code>pg_database .oid</code>)	OID of database in which the statement was executed
<code>toplevel bool</code>	True if the query was executed as a top-level statement (always true if <code>pg_stat_statements.track</code> is set to <code>top</code>)
<code>queryid bigint</code>	Hash code to identify identical normalized queries.
<code>query text</code>	Text of a representative statement
<code>plans bigint</code>	Number of times the statement was planned (if <code>pg_stat_statements.track_planning</code> is enabled, otherwise zero)
<code>total_plan_time double precision</code>	Total time spent planning the statement, in milliseconds (if <code>pg_stat_statements.track_planning</code> is enabled, otherwise zero)
<code>min_plan_time double precision</code>	Minimum time spent planning the statement, in milliseconds. This field will be zero if <code>pg_stat_statements.track_planning</code> is disabled, or if the counter has been reset using the <code>pg_stat_statements_reset</code> function with the <code>minmax_only</code> parameter set to <code>true</code> and never been planned since.
<code>max_plan_time double precision</code>	Maximum time spent planning the statement, in milliseconds. This field will be zero if <code>pg_stat_statements.track_planning</code> is disabled, or if the counter has been reset using the

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Column Type	Description
<code>pg_stat_statements_reset</code>	function with the <code>minmax_only</code> parameter set to <code>true</code> and never been planned since.
<code>mean_plan_time</code> double precision	Mean time spent planning the statement, in milliseconds (if <code>pg_stat_statements.track_planning</code> is enabled, otherwise zero)
<code>stddev_plan_time</code> double precision	Population standard deviation of time spent planning the statement, in milliseconds (if <code>pg_stat_statements.track_planning</code> is enabled, otherwise zero)
<code>calls</code> bigint	Number of times the statement was executed
<code>total_exec_time</code> double precision	Total time spent executing the statement, in milliseconds
<code>min_exec_time</code> double precision	Minimum time spent executing the statement, in milliseconds, this field will be zero until this statement is executed first time after reset performed by the <code>pg_stat_statements_reset</code> function with the <code>minmax_only</code> parameter set to <code>true</code>
<code>max_exec_time</code> double precision	Maximum time spent executing the statement, in milliseconds, this field will be zero until this statement is executed first time after reset performed by the <code>pg_stat_statements_reset</code> function with the <code>minmax_only</code> parameter set to <code>true</code>
<code>mean_exec_time</code> double precision	Mean time spent executing the statement, in milliseconds
<code>stddev_exec_time</code> double precision	Population standard deviation of time spent executing the statement, in milliseconds
<code>rows</code> bigint	Total number of rows retrieved or affected by the statement
<code>shared_blks_hit</code> bigint	Total number of shared block cache hits by the statement
<code>shared_blks_read</code> bigint	Total number of shared blocks read by the statement
<code>shared_blks_dirtied</code> bigint	Total number of shared blocks dirtied by the statement
<code>shared_blks_written</code> bigint	Total number of shared blocks written by the statement
<code>local_blks_hit</code> bigint	Total number of local block cache hits by the statement
<code>local_blks_read</code> bigint	Total number of local blocks read by the statement
<code>local_blks_dirtied</code> bigint	Total number of local blocks dirtied by the statement
<code>local_blks_written</code> bigint	Total number of local blocks written by the statement
<code>temp_blks_read</code> bigint	Total number of temp blocks read by the statement
<code>temp_blks_written</code> bigint	Total number of temp blocks written by the statement
<code>shared_blk_read_time</code> double precision	

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

Column Type	Description
	Total time the statement spent reading shared blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>shared_blk_write_time</code> double precision	Total time the statement spent writing shared blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>local_blk_read_time</code> double precision	Total time the statement spent reading local blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>local_blk_write_time</code> double precision	Total time the statement spent writing local blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>temp_blk_read_time</code> double precision	Total time the statement spent reading temporary file blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>temp_blk_write_time</code> double precision	Total time the statement spent writing temporary file blocks, in milliseconds (if <code>track_io_timing</code> is enabled, otherwise zero)
<code>wal_records</code> bigint	Total number of WAL records generated by the statement
<code>wal_fpi</code> bigint	Total number of WAL full page images generated by the statement
<code>wal_bytes</code> numeric	Total amount of WAL generated by the statement in bytes
<code>jit_functions</code> bigint	Total number of functions JIT-compiled by the statement
<code>jit_generation_time</code> double precision	Total time spent by the statement on generating JIT code, in milliseconds
<code>jit_inlining_count</code> bigint	Number of times functions have been inlined
<code>jit_inlining_time</code> double precision	Total time spent by the statement on inlining functions, in milliseconds
<code>jit_optimization_count</code> bigint	Number of times the statement has been optimized
<code>jit_optimization_time</code> double precision	Total time spent by the statement on optimizing, in milliseconds
<code>jit_emission_count</code> bigint	Number of times code has been emitted
<code>jit_emission_time</code> double precision	Total time spent by the statement on emitting code, in milliseconds
<code>jit_deform_count</code> bigint	Total number of tuple deform functions JIT-compiled by the statement
<code>jit_deform_time</code> double precision	Total time spent by the statement on JIT-compiling tuple deform functions, in milliseconds
<code>stats_since</code> timestamp with time zone	Time at which statistics gathering started for this statement
<code>minmax_stats_since</code> timestamp with time zone	

Column Type	Description
	Time at which min/max statistics gathering started for this statement (fields <code>min_plan_time</code> , <code>max_plan_time</code> , <code>min_exec_time</code> and <code>max_exec_time</code>)

For security reasons, only superusers and roles with privileges of the `pg_read_all_stats` role are allowed to see the SQL text and `queryid` of queries executed by other users. Other users can see the statistics, however, if the view has been installed in their database.

Plannable queries (that is, `SELECT`, `INSERT`, `UPDATE`, `DELETE`, and `MERGE`) and utility commands are combined into a single `pg_stat_statements` entry whenever they have identical query structures according to an internal hash calculation. Typically, two queries will be considered the same for this purpose if they are semantically equivalent except for the values of literal constants appearing in the query.

Note

The following details about constant replacement and `queryid` only apply when `compute_query_id` is enabled. If you use an external module instead to compute `queryid`, you should refer to its documentation for details.

When a constant's value has been ignored for purposes of matching the query to other queries, the constant is replaced by a parameter symbol, such as `$1`, in the `pg_stat_statements` display. The rest of the query text is that of the first query that had the particular `queryid` hash value associated with the `pg_stat_statements` entry.

Queries on which normalization can be applied may be observed with constant values in `pg_stat_statements`, especially when there is a high rate of entry deallocations. To reduce the likelihood of this happening, consider increasing `pg_stat_statements.max`. The `pg_stat_statements_info` view, discussed below in [Section F.9.2](#), provides statistics about entry deallocations.

In some cases, queries with visibly different texts might get merged into a single `pg_stat_statements` entry. Normally this will happen only for semantically equivalent queries, but there is a small chance of hash collisions causing unrelated queries to be merged into one entry. (This cannot happen for queries belonging to different users or databases, however.)

Since the `queryid` hash value is computed on the post-parse-analysis representation of the queries, the opposite is also possible: queries with identical texts might appear as separate entries, if they have different meanings as a result of factors such as different `search_path` settings.

Consumers of `pg_stat_statements` may wish to use `queryid` (perhaps in combination with `dbid` and `userid`) as a more stable and reliable identifier for each entry than its query text. However, it is important to understand that there are only limited guarantees around the stability of the `queryid` hash value. Since the identifier is derived from the post-parse-analysis tree, its value is a function of, among other things, the internal object identifiers appearing in this representation. This has some counterintuitive implications. For example, `pg_stat_statements` will consider two apparently-identical queries to be distinct, if they reference a table that was dropped and recreated between the executions of the two queries. The hashing process is also sensitive to differences in machine architecture and other facets of the platform. Furthermore, it is not safe to assume that `queryid` will be stable across major versions of Postgres Pro.

Two servers participating in replication based on physical WAL replay can be expected to have identical `queryid` values for the same query. However, logical replication schemes do not promise to keep replicas identical in all relevant details, so `queryid` will not be a useful identifier for accumulating costs across a set of logical replicas. If in doubt, direct testing is recommended.

Generally, it can be assumed that `queryid` values are stable between minor version releases of Postgres Pro, providing that instances are running on the same machine architecture and the catalog metadata details match. Compatibility will only be broken between minor versions as a last resort.

The parameter symbols used to replace constants in representative query texts start from the next number after the highest `$n` parameter in the original query text, or `$1` if there was none. It's worth noting that in some cases there may be hidden parameter symbols that affect this numbering. For example, PL/pgSQL uses hidden parameter symbols to insert values of function local variables into queries, so that a PL/pgSQL statement like `SELECT i + 1 INTO j` would have representative text like `SELECT i + $2`.

The representative query texts are kept in an external disk file, and do not consume shared memory. Therefore, even very lengthy query texts can be stored successfully. However, if many long query texts are accumulated, the external file might grow unmanageably large. As a recovery method if that happens, `pg_stat_statements` may choose to discard the query texts, whereupon all existing entries in the `pg_stat_statements` view will show null `query` fields, though the statistics associated with each `queryid` are preserved. If this happens, consider reducing `pg_stat_statements.max` to prevent recurrences.

`plans` and `calls` aren't always expected to match because planning and execution statistics are updated at their respective end phase, and only for successful operations. For example, if a statement is successfully planned but fails during the execution phase, only its planning statistics will be updated. If planning is skipped because a cached plan is used, only its execution statistics will be updated.

F.9.2. The `pg_stat_statements_info` View

The statistics of the `pg_stat_statements` module itself are tracked and made available via a view named `pg_stat_statements_info`. This view contains only a single row. The columns of the view are shown in [Table F.24](#).

Table F.24. `pg_stat_statements_info` Columns

Column Type	Description
<code>dealloc bigint</code>	Total number of times <code>pg_stat_statements</code> entries about the least-executed statements were deallocated because more distinct statements than <code>pg_stat_statements.max</code> were observed
<code>stats_reset timestamp with time zone</code>	Time at which all statistics in the <code>pg_stat_statements</code> view were last reset.

F.9.3. Functions

`pg_stat_statements_reset(userid Oid, dbid Oid, queryid bigint, minmax_only boolean)` returns `timestamp with time zone`

`pg_stat_statements_reset` discards statistics gathered so far by `pg_stat_statements` corresponding to the specified `userid`, `dbid` and `queryid`. If any of the parameters are not specified, the default value `0`(invalid) is used for each of them and the statistics that match with other parameters will be reset. If no parameter is specified or all the specified parameters are `0`(invalid), it will discard all statistics. If all statistics in the `pg_stat_statements` view are discarded, it will also reset the statistics in the `pg_stat_statements_info` view. When `minmax_only` is `true` only the values of minimum and maximum planning and execution time will be reset (i.e. `min_plan_time`, `max_plan_time`, `min_exec_time` and `max_exec_time` fields). The default value for `minmax_only` parameter is `false`. Time of last min/max reset performed is shown in `minmax_stats_since` field of the `pg_stat_statements` view. This function returns the time of a reset. This time is saved to `stats_reset` field of `pg_stat_statements_info` view or to `minmax_stats_since` field of the `pg_stat_statements` view if the corresponding reset was actually performed. By default, this function can only be executed by superusers. Access may be granted to others using `GRANT`.

`pg_stat_statements(showtext boolean)` returns `setof record`

The `pg_stat_statements` view is defined in terms of a function also named `pg_stat_statements`. It is possible for clients to call the `pg_stat_statements` function directly, and by specifying `showtext := false` have query text be omitted (that is, the `OUT` argument that corresponds to the view's `query` column will return nulls). This feature is intended to support external tools that might wish to avoid

the overhead of repeatedly retrieving query texts of indeterminate length. Such tools can instead cache the first query text observed for each entry themselves, since that is all `pg_stat_statements` itself does, and then retrieve query texts only as needed. Since the server stores query texts in a file, this approach may reduce physical I/O for repeated examination of the `pg_stat_statements` data.

F.9.4. Configuration Parameters

`pg_stat_statements.max` (integer)

`pg_stat_statements.max` is the maximum number of statements tracked by the module (i.e., the maximum number of rows in the `pg_stat_statements` view). If more distinct statements than that are observed, information about the least-executed statements is discarded. The number of times such information was discarded can be seen in the `pg_stat_statements_info` view. The default value is 5000. This parameter can only be set at server start.

`pg_stat_statements.track` (enum)

`pg_stat_statements.track` controls which statements are counted by the module. Specify `top` to track top-level statements (those issued directly by clients), `all` to also track nested statements (such as statements invoked within functions), or `none` to disable statement statistics collection. The default value is `top`. Only superusers can change this setting.

`pg_stat_statements.track_utility` (boolean)

`pg_stat_statements.track_utility` controls whether utility commands are tracked by the module. Utility commands are all those other than `SELECT`, `INSERT`, `UPDATE`, `DELETE`, and `MERGE`. The default value is `on`. Only superusers can change this setting.

`pg_stat_statements.track_planning` (boolean)

`pg_stat_statements.track_planning` controls whether planning operations and duration are tracked by the module. Enabling this parameter may incur a noticeable performance penalty, especially when statements with identical query structure are executed by many concurrent connections which compete to update a small number of `pg_stat_statements` entries. The default value is `off`. Only superusers can change this setting.

`pg_stat_statements.save` (boolean)

`pg_stat_statements.save` specifies whether to save statement statistics across server shutdowns. If it is `off` then statistics are not saved at shutdown nor reloaded at server start. The default value is `on`. This parameter can only be set in the `postgresql.conf` file or on the server command line.

The module requires additional shared memory proportional to `pg_stat_statements.max`. Note that this memory is consumed whenever the module is loaded, even if `pg_stat_statements.track` is set to `none`.

These parameters must be set in `postgresql.conf`. Typical usage might be:

```
# postgresql.conf
shared_preload_libraries = 'pg_stat_statements'

compute_query_id = on
pg_stat_statements.max = 10000
pg_stat_statements.track = all
```

F.9.5. Sample Output

```
bench=# SELECT pg_stat_statements_reset();

$ pgbench -i bench
$ pgbench -c10 -t300 bench

bench=# \x
bench=# SELECT query, calls, total_exec_time, rows, 100.0 * shared_blks_hit /
        nullif(shared_blks_hit + shared_blks_read, 0) AS hit_percent
```

Additional Supplied Modules and Extensions Shipped in

postgrespro-sdm-17-contrib

FROM pg_stat_statements ORDER BY total_exec_time DESC LIMIT 5;

```

-[ RECORD 1 ]-----+-----
query          | UPDATE pgbench_branches SET bbalance = bbalance + $1 WHERE bid = $2
calls          | 3000
total_exec_time | 25565.855387
rows          | 3000
hit_percent    | 100.0000000000000000
-[ RECORD 2 ]-----+-----
query          | UPDATE pgbench_tellers SET tbalance = tbalance + $1 WHERE tid = $2
calls          | 3000
total_exec_time | 20756.669379
rows          | 3000
hit_percent    | 100.0000000000000000
-[ RECORD 3 ]-----+-----
query          | copy pgbench_accounts from stdin
calls          | 1
total_exec_time | 291.865911
rows          | 100000
hit_percent    | 100.0000000000000000
-[ RECORD 4 ]-----+-----
query          | UPDATE pgbench_accounts SET abalance = abalance + $1 WHERE aid = $2
calls          | 3000
total_exec_time | 271.232977
rows          | 3000
hit_percent    | 98.8454011741682975
-[ RECORD 5 ]-----+-----
query          | alter table pgbench_accounts add primary key (aid)
calls          | 1
total_exec_time | 160.588563
rows          | 0
hit_percent    | 100.0000000000000000

```

```

bench=# SELECT pg_stat_statements_reset(0,0,s.queryid) FROM pg_stat_statements AS s
        WHERE s.query = 'UPDATE pgbench_branches SET bbalance = bbalance + $1 WHERE
        bid = $2';

```

```

bench=# SELECT query, calls, total_exec_time, rows, 100.0 * shared_blks_hit /
        nullif(shared_blks_hit + shared_blks_read, 0) AS hit_percent
        FROM pg_stat_statements ORDER BY total_exec_time DESC LIMIT 5;

```

```

-[ RECORD 1 ]-----+-----
query          | UPDATE pgbench_tellers SET tbalance = tbalance + $1 WHERE tid = $2
calls          | 3000
total_exec_time | 20756.669379
rows          | 3000
hit_percent    | 100.0000000000000000
-[ RECORD 2 ]-----+-----
query          | copy pgbench_accounts from stdin
calls          | 1
total_exec_time | 291.865911
rows          | 100000
hit_percent    | 100.0000000000000000
-[ RECORD 3 ]-----+-----
query          | UPDATE pgbench_accounts SET abalance = abalance + $1 WHERE aid = $2
calls          | 3000
total_exec_time | 271.232977
rows          | 3000
hit_percent    | 98.8454011741682975

```

Additional Supplied Modules
and Extensions Shipped in
postgrespro-sdm-17-contrib

```

-[ RECORD 4 ]---+-----
query          | alter table pgbench_accounts add primary key (aid)
calls           | 1
total_exec_time | 160.588563
rows           | 0
hit_percent     | 100.0000000000000000
-[ RECORD 5 ]---+-----
query          | vacuum analyze pgbench_accounts
calls           | 1
total_exec_time | 136.448116
rows           | 0
hit_percent     | 99.9201915403032721

bench=# SELECT pg_stat_statements_reset(0,0,0);

bench=# SELECT query, calls, total_exec_time, rows, 100.0 * shared_blks_hit /
        nullif(shared_blks_hit + shared_blks_read, 0) AS hit_percent
        FROM pg_stat_statements ORDER BY total_exec_time DESC LIMIT 5;
-[ RECORD 1 ]---+-----
query          | SELECT pg_stat_statements_reset(0,0,0)
calls           | 1
total_exec_time | 0.189497
rows           | 1
hit_percent     | 
-[ RECORD 2 ]---+-----
query          | SELECT query, calls, total_exec_time, rows, $1 * shared_blks_hit /
+              | nullif(shared_blks_hit + shared_blks_read, $2) AS
hit_percent+   | FROM pg_stat_statements ORDER BY total_exec_time DESC LIMIT
$3
calls           | 0
total_exec_time | 0
rows           | 0
hit_percent     | 

```

F.9.6. Authors

Takahiro Itagaki <itagaki.takahiro@oss.ntt.co.jp>. Query normalization added by Peter Geoghegan <peter@2ndquadrant.com>.

F.10. pg_wait_sampling — collecting sampling-based statistics on wait events

`pg_wait_sampling` is a Postgres Pro extension for collecting sampling-based statistics on wait events.

Starting from the 9.6 version, Postgres Pro Shardman provides information about the current wait events for particular processes. However, to get descriptive statistics of the server activity, you have to sample wait events multiple times. The `pg_wait_sampling` extension automates wait events sampling by launching a special background worker. With `pg_wait_sampling` enabled, you can get the following sampling-based data:

- *Waits history* — the list of wait events for recent processes, with timestamps.
- *Waits profile* — the number of wait event samples for all processes over time (per wait event type).
- Current wait events for all processes, including background workers.

Using `pg_wait_sampling`, you can troubleshoot dependencies for queries that process longer than expected. You can see what a particular process is waiting for at each moment of time, and analyze wait events statistics. For the list of possible wait events, see [Table 27.4](#).

In combination with `pg_stat_statements`, this extension can also provide per query statistics.

See Also

[Viewing Statistics on Wait Events](#)

[Reference](#)

F.10.1. Installation

The `pg_wait_sampling` extension is included into Postgres Pro Shardman and requires no special prerequisites.

To enable `pg_wait_sampling`, do the following:

1. Add `pg_wait_sampling` to the `shared_preload_libraries` variable in the `postgresql.conf` file:

```
shared_preload_libraries = 'pg_wait_sampling'
```

Important

When using the `pg_wait_sampling` with the `pg_stat_statements` extension, place `pg_stat_statements` before `pg_wait_sampling` in the `shared_preload_libraries` variable. This ensures that the `pg_stat_statements` extension does not overwrite query IDs that are used by `pg_wait_sampling`.

2. Create the `pg_wait_sampling` extension using the following query:

```
CREATE EXTENSION pg_wait_sampling;
```

3. Restart the server. Since `pg_wait_sampling` requires additional shared memory and launches a background worker, you must restart the server after adding or removing `pg_wait_sampling`.

Once the extension is created, `pg_wait_sampling` starts collecting statistics on wait events for each process. Optionally, you can change the sampling frequency and configure statistics collection modes using [GUC variables](#).

If you want to stop collecting statistics, drop the `pg_wait_sampling` extension, remove the `pg_wait_sampling` from the `shared_preload_libraries` variable, and restart the server again.

F.10.2. Usage

F.10.2.1. Viewing Statistics on Wait Events

When `pg_wait_sampling` is enabled, wait events are sampled each 10ms by default. You can access the collected statistics through a set of functions and views. If some of your queries appear stalled or take longer than expected, you can review all the current wait events in the `pg_wait_sampling_current` view:

```
SELECT * FROM pg_wait_sampling_current;
```

The returned statistics covers all the active processes, including background workers. To filter the view for a single process, run `pg_wait_sampling_get_current()` function specifying the process ID:

```
SELECT * FROM pg_wait_sampling_get_current(pid);
```

To better understand the dynamics of the recent wait events, you can access the `pg_wait_sampling_history` view that provides the wait events history for the latest processes:

```
SELECT * FROM pg_wait_sampling_history;
```

The returned view displays wait samples per process, with event timestamps. Waits history is implemented as an in-memory ring buffer. By default, the history size is limited to 5000 samples. To save history for a longer period of time, you can modify the `pg_wait_sampling.history_size` variable, or connect a client application to periodically read waits history and dump it to a local storage.

To monitor wait events in your database over time, use the `pg_wait_sampling_profile` view:

```
SELECT * FROM pg_wait_sampling_profile;
```

Waits profile is stored as an in-memory hash table that accumulates samples per each process and each wait event. You can reset waits profile by calling the `pg_wait_sampling_reset_profile()` function. If you connect a client to your database that periodically dumps the wait events data and resets the profile, you can save and compare statistics of wait events intensity over time.

Important

Since wait sampling statistics is stored in RAM, make sure to reset the waits profile from time to time to avoid memory overflow.

If you are not interested in the distribution of wait events between processes, you can set the `pg_wait_sampling.profile_pid` variable to `false`. In this case, the `pid` value for each process is set to zero, and wait samples for all the processes are stored together.

F.10.3. Reference

F.10.3.1. GUC Variables

The `pg_wait_sampling` extension is configured with GUC variables. They are stored in the local memory and can be changed in the `postgresql.conf` file. For the changes to take effect, reload the server configuration. You can do it, for example, by calling the `pg_reload_conf()` function or executing the `pg_ctl reload` command.

`pg_wait_sampling.history_size (int4)`

Size of in-memory ring buffer for history sampling, in the number of samples.

Default: 5000

`pg_wait_sampling.history_period (int4)`

Time interval for history sampling, in milliseconds.

Default: 10

`pg_wait_sampling.profile_period (int4)`

Time interval for waits profile sampling, in milliseconds.

Default: 10

`pg_wait_sampling.profile_pid (bool)`

Sampling mode for waits profile. If set to `true`, waits profile is collected per process. If you set `pg_wait_sampling.profile_pid` to `false`, waits profile is collected for all processes together. In this case, the `pid` value for each process is set to zero, and the corresponding row contains wait samples for all the processes.

Default: `true`

`pg_wait_sampling.profile_queries (enum)`

Sampling mode for the waits profile. If `pg_wait_sampling.profile_queries` is set to `none`, the `queryid` field in views will be zero. If it is set to `top`, query IDs of only the top level statements are recorded. If it is set to `all`, query IDs of nested statements are recorded. To collect the waits profile per query, ensure that the [pg_stat_statements](#) extension is configured and set `pg_wait_sampling.profile_queries` to `top`. For version 14 or higher, you can compute query IDs by configuring the `compute_query_id` parameter and set the `pg_wait_sampling.profile_queries` value to `top` or `all`. For details, see [Section 19.9.2](#).

Default: `top`

`pg_wait_sampling.sample_cpu (bool)`

The sampling mode that determines whether to perform sampling of on-CPU backends. If `pg_wait_sampling.sample_cpu` is set to `true`, then sampling also includes processes that are not waiting for anything. The wait event columns for such processes will have a `NULL` value.

Default: `true`

F.10.3.2. pg_wait_sampling Views

F.10.3.2.1. pg_wait_sampling_current View

The `pg_wait_sampling_current` view provides the information about the current wait events for all processes, including background workers.

Table F.25. pg_wait_sampling_current View

Column Name	Column Type	Description
<code>pid</code>	<code>int4</code>	Process ID
<code>event_type</code>	<code>text</code>	Name of wait event type
<code>event</code>	<code>text</code>	Name of wait event
<code>queryid</code>	<code>int8</code>	Query ID

F.10.3.2.2. pg_wait_sampling_history View

The `pg_wait_sampling_history` view provides the history of wait events. This data is stored as an in-memory ring buffer.

Table F.26. pg_wait_sampling_history View

Column Name	Column Type	Description
<code>pid</code>	<code>int4</code>	Process ID

Column Name	Column Type	Description
ts	timestampz	Sample timestamp
event_type	text	Name of wait event type
event	text	Name of wait event
queryid	int8	Query ID

F.10.3.2.3. `pg_wait_sampling_profile` View

The `pg_wait_sampling_profile` view provides the profile of wait events. This data is stored as an in-memory hash table.

Table F.27. `pg_wait_sampling_profile` View

Column Name	Column Type	Description
pid	int4	Process ID
event_type	text	Name of wait event type
event	text	Name of wait event
queryid	int8	Query ID
count	int8	Number of samples

F.10.3.3. Functions

`pg_wait_sampling_get_current(pid int4)`

Returns the `pg_wait_sampling_current` view with the list of current wait events. If you set the *pid* argument, the view is filtered for the process with this *pid*.

Arguments:

- *pid* — Optional. The process ID for which to display the current wait events.

`pg_wait_sampling_reset_profile()`

Resets the waits profile and clears the memory.

F.10.4. Authors

- Alexander Korotkov
- Ildus Kurbangaliev <i.kurbangaliev@gmail.com>

F.11. postgres_fdw — access data stored in external Postgres Pro servers

The `postgres_fdw` module provides the foreign-data wrapper `postgres_fdw`, which can be used to access data stored in external Postgres Pro servers.

Deleting this extension with `CASCADE` will result in Postgres Pro Shardman failure.

The functionality provided by this module overlaps substantially with the functionality of the older [dblink](#) module. But `postgres_fdw` provides more transparent and standards-compliant syntax for accessing remote tables, and can give better performance in many cases.

Note

Note that you cannot `ALTER` or `DROP` automatically created foreign objects. However, you can `CREATE`, `ALTER`, and `DELETE` newly created foreign objects.

To prepare for remote access using `postgres_fdw`:

1. Install the `postgres_fdw` extension using [CREATE EXTENSION](#).
2. Create a foreign server object, using [CREATE SERVER](#), to represent each remote database you want to connect to. Specify connection information, except `user` and `password`, as options of the server object.
3. Create a user mapping with `CREATE USER MAPPING` for each database user you want to allow to access each foreign server. Specify the remote user name and password to use as `user` and `password` options of the user mapping.
4. Create a foreign table with `CREATE FOREIGN TABLE` for each remote table you want to access. The columns of the foreign table must match the referenced remote table. You can, however, use table and/or column names different from the remote table's, if you specify the correct remote names as options of the foreign table object.

Now you need only `SELECT` from a foreign table to access the data stored in its underlying remote table. You can also modify the remote table using `INSERT`, `UPDATE`, `DELETE`, `COPY`, or `TRUNCATE`. (Of course, the remote user you have specified in your user mapping must have privileges to do these things.)

Note that the `ONLY` option specified in `SELECT`, `UPDATE`, `DELETE` or `TRUNCATE` has no effect when accessing or modifying the remote table.

Note that `postgres_fdw` currently lacks support for `INSERT` statements with an `ON CONFLICT DO UPDATE` clause. However, the `ON CONFLICT DO NOTHING` clause is supported, provided a unique index inference specification is omitted. Note also that `postgres_fdw` supports row movement invoked by `UPDATE` statements executed on partitioned tables, but it currently does not handle the case where a remote partition chosen to insert a moved row into is also an `UPDATE` target partition that will be updated elsewhere in the same command.

It is generally recommended that the columns of a foreign table be declared with exactly the same data types, and collations if applicable, as the referenced columns of the remote table. Although `postgres_fdw` is currently rather forgiving about performing data type conversions at need, surprising semantic anomalies may arise when types or collations do not match, due to the remote server interpreting query conditions differently from the local server.

Note that a foreign table can be declared with fewer columns, or with a different column order, than its underlying remote table has. Matching of columns to the remote table is by name, not position.

F.11.1. FDW Options of `postgres_fdw`

F.11.1.1. Connection Options

A foreign server using the `postgres_fdw` foreign data wrapper can have the same options that `libpq` accepts in connection strings, as described in [Section 35.1.2](#), except that these options are not allowed or have special handling:

- `user`, `password` and `sslpassword` (specify these in a user mapping instead, or use a service file)
- `client_encoding` (this is automatically set from the local server encoding)
- `application_name` - this may appear in *either or both* a connection and [postgres_fdw.application_name](#). If both are present, `postgres_fdw.application_name` overrides the connection setting. Unlike `libpq`, `postgres_fdw` allows `application_name` to include “escape sequences”. See [postgres_fdw.application_name](#) for details.
- `fallback_application_name` (always set to `postgres_fdw`)
- `sslkey` and `sslcert` - these may appear in *either or both* a connection and a user mapping. If both are present, the user mapping setting overrides the connection setting.

Only superusers may create or modify user mappings with the `sslcert` or `sslkey` settings.

Non-superusers may connect to foreign servers using password authentication or with GSSAPI delegated credentials, so specify the `password` option for user mappings belonging to non-superusers where password authentication is required.

A superuser may override this check on a per-user-mapping basis by setting the user mapping option `password_required 'false'`, e.g.,

```
ALTER USER MAPPING FOR some_non_superuser SERVER loopback_nopw
OPTIONS (ADD password_required 'false');
```

To prevent unprivileged users from exploiting the authentication rights of the unix user the `postgres` server is running as to escalate to superuser rights, only the superuser may set this option on a user mapping.

Care is required to ensure that this does not allow the mapped user the ability to connect as superuser to the mapped database per CVE-2007-3278 and CVE-2007-6601. Don't set `password_required=false` on the `public` role. Keep in mind that the mapped user can potentially use any client certificates, `.pgpass`, `.pg_service.conf` etc. in the unix home directory of the system user the `postgres` server runs as. They can also use any trust relationship granted by authentication modes like `peer` or `ident` authentication.

F.11.1.2. Object Name Options

These options can be used to control the names used in SQL statements sent to the remote Postgres Pro server. These options are needed when a foreign table is created with names different from the underlying remote table's names.

`schema_name` (string)

This option, which can be specified for a foreign table, gives the schema name to use for the foreign table on the remote server. If this option is omitted, the name of the foreign table's schema is used.

`table_name` (string)

This option, which can be specified for a foreign table, gives the table name to use for the foreign table on the remote server. If this option is omitted, the foreign table's name is used.

`column_name` (string)

This option, which can be specified for a column of a foreign table, gives the column name to use for the column on the remote server. If this option is omitted, the column's name is used.

F.11.1.3. Cost Estimation Options

`postgres_fdw` retrieves remote data by executing queries against remote servers, so ideally the estimated cost of scanning a foreign table should be whatever it costs to be done on the remote server, plus

some overhead for communication. The most reliable way to get such an estimate is to ask the remote server and then add something for overhead — but for simple queries, it may not be worth the cost of an additional remote query to get a cost estimate. So `postgres_fdw` provides the following options to control how cost estimation is done:

`use_remote_estimate` (boolean)

This option, which can be specified for a foreign table or a foreign server, controls whether `postgres_fdw` issues remote `EXPLAIN` commands to obtain cost estimates. A setting for a foreign table overrides any setting for its server, but only for that table. The default is `false`.

`fdw_startup_cost` (floating point)

This option, which can be specified for a foreign server, is a floating point value that is added to the estimated startup cost of any foreign-table scan on that server. This represents the additional overhead of establishing a connection, parsing and planning the query on the remote side, etc. The default value is 100.

`fdw_tuple_cost` (floating point)

This option, which can be specified for a foreign server, is a floating point value that is used as extra cost per-tuple for foreign-table scans on that server. This represents the additional overhead of data transfer between servers. You might increase or decrease this number to reflect higher or lower network delay to the remote server. The default value is 0.2.

When `use_remote_estimate` is `true`, `postgres_fdw` obtains row count and cost estimates from the remote server and then adds `fdw_startup_cost` and `fdw_tuple_cost` to the cost estimates. When `use_remote_estimate` is `false`, `postgres_fdw` performs local row count and cost estimation and then adds `fdw_startup_cost` and `fdw_tuple_cost` to the cost estimates. This local estimation is unlikely to be very accurate unless local copies of the remote table's statistics are available. Running [ANALYZE](#) on the foreign table is the way to update the local statistics; this will perform a scan of the remote table and then calculate and store statistics just as though the table were local. Keeping local statistics can be a useful way to reduce per-query planning overhead for a remote table — but if the remote table is frequently updated, the local statistics will soon be obsolete.

The following option controls how such an `ANALYZE` operation behaves:

`analyze_sampling` (string)

This option, which can be specified for a foreign table or a foreign server, determines if `ANALYZE` on a foreign table samples the data on the remote side, or reads and transfers all data and performs the sampling locally. The supported values are `off`, `random`, `system`, `bernoulli` and `auto`. `off` disables remote sampling, so all data are transferred and sampled locally. `random` performs remote sampling using the `random()` function to choose returned rows, while `system` and `bernoulli` rely on the built-in `TABLESAMPLE` methods of those names. `random` works on all remote server versions, while `TABLESAMPLE` is supported only since 9.5. `auto` (the default) picks the recommended sampling method automatically; currently it means either `bernoulli` or `random` depending on the remote server version.

F.11.1.4. Remote Execution Options

By default, only `WHERE` clauses using built-in operators and functions will be considered for execution on the remote server. Clauses involving non-built-in functions are checked locally after rows are fetched. If such functions are available on the remote server and can be relied on to produce the same results as they do locally, performance can be improved by sending such `WHERE` clauses for remote execution. This behavior can be controlled using the following option:

`extensions` (string)

This option is a comma-separated list of names of Postgres Pro extensions that are installed, in compatible versions, on both the local and remote servers. Functions and operators that are immutable and belong to a listed extension will be considered shippable to the remote server. This option can only be specified for foreign servers, not per-table.

When using the `extensions` option, *it is the user's responsibility* that the listed extensions exist and behave identically on both the local and remote servers. Otherwise, remote queries may fail or behave unexpectedly.

`fetch_size` (integer)

This option specifies the number of rows `postgres_fdw` should get in each fetch operation. It can be specified for a foreign table or a foreign server. The option specified on a table overrides an option specified for the server. The default is 100.

`batch_size` (integer)

This option specifies the number of rows `postgres_fdw` should insert in each insert operation. It can be specified for a foreign table or a foreign server. The option specified on a table overrides an option specified for the server. The default is 1.

Note the actual number of rows `postgres_fdw` inserts at once depends on the number of columns and the provided `batch_size` value. The batch is executed as a single query, and the libpq protocol (which `postgres_fdw` uses to connect to a remote server) limits the number of parameters in a single query to 65535. When the number of columns * `batch_size` exceeds the limit, the `batch_size` will be adjusted to avoid an error.

This option also applies when copying into foreign tables. In that case the actual number of rows `postgres_fdw` copies at once is determined in a similar way to the insert case, but it is limited to at most 1000 due to implementation restrictions of the `COPY` command.

F.11.1.5. Asynchronous Execution Options

`postgres_fdw` supports asynchronous execution, which runs multiple parts of an `Append` node concurrently rather than serially to improve performance. This execution can be controlled using the following option:

`async_capable` (boolean)

This option controls whether `postgres_fdw` allows foreign tables to be scanned concurrently for asynchronous execution. It can be specified for a foreign table or a foreign server. A table-level option overrides a server-level option. The default is `false`.

In order to ensure that the data being returned from a foreign server is consistent, `postgres_fdw` will only open one connection for a given foreign server and will run all queries against that server sequentially even if there are multiple foreign tables involved, unless those tables are subject to different user mappings. In such a case, it may be more performant to disable this option to eliminate the overhead associated with running queries asynchronously.

Asynchronous execution is applied even when an `Append` node contains subplan(s) executed synchronously as well as subplan(s) executed asynchronously. In such a case, if the asynchronous subplans are ones processed using `postgres_fdw`, tuples from the asynchronous subplans are not returned until after at least one synchronous subplan returns all tuples, as that subplan is executed while the asynchronous subplans are waiting for the results of asynchronous queries sent to foreign servers. This behavior might change in a future release.

F.11.1.6. Transaction Management Options

As described in the Transaction Management section, in `postgres_fdw` transactions are managed by creating corresponding remote transactions, and subtransactions are managed by creating corresponding remote subtransactions. When multiple remote transactions are involved in the current local transaction, by default `postgres_fdw` commits or aborts those remote transactions serially when the local transaction is committed or aborted. When multiple remote subtransactions are involved in the current local subtransaction, by default `postgres_fdw` commits or aborts those remote subtransactions serially when the local subtransaction is committed or aborted. Performance can be improved with the following options:

`parallel_commit` (boolean)

This option controls whether `postgres_fdw` commits, in parallel, remote transactions opened on a foreign server in a local transaction when the local transaction is committed. This setting also applies to remote and local subtransactions. This option can only be specified for foreign servers, not per-table. The default is `false`.

`parallel_abort` (boolean)

This option controls whether `postgres_fdw` aborts, in parallel, remote transactions opened on a foreign server in a local transaction when the local transaction is aborted. This setting also applies to remote and local subtransactions. This option can only be specified for foreign servers, not per-table. The default is `false`.

If multiple foreign servers with these options enabled are involved in a local transaction, multiple remote transactions on those foreign servers are committed or aborted in parallel across those foreign servers when the local transaction is committed or aborted.

When these options are enabled, a foreign server with many remote transactions may see a negative performance impact when the local transaction is committed or aborted.

F.11.1.7. Updatability Options

By default all foreign tables using `postgres_fdw` are assumed to be updatable. This may be overridden using the following option:

`updatable` (boolean)

This option controls whether `postgres_fdw` allows foreign tables to be modified using `INSERT`, `UPDATE` and `DELETE` commands. It can be specified for a foreign table or a foreign server. A table-level option overrides a server-level option. The default is `true`.

Of course, if the remote table is not in fact updatable, an error would occur anyway. Use of this option primarily allows the error to be thrown locally without querying the remote server. Note however that the `information_schema` views will report a `postgres_fdw` foreign table to be updatable (or not) according to the setting of this option, without any check of the remote server.

F.11.1.8. Truncatability Options

By default all foreign tables using `postgres_fdw` are assumed to be truncatable. This may be overridden using the following option:

`truncatable` (boolean)

This option controls whether `postgres_fdw` allows foreign tables to be truncated using the `TRUNCATE` command. It can be specified for a foreign table or a foreign server. A table-level option overrides a server-level option. The default is `true`.

Of course, if the remote table is not in fact truncatable, an error would occur anyway. Use of this option primarily allows the error to be thrown locally without querying the remote server.

F.11.1.9. Importing Options

`postgres_fdw` is able to import foreign table definitions using `IMPORT FOREIGN SCHEMA`. This command creates foreign table definitions on the local server that match tables or views present on the remote server. If the remote tables to be imported have columns of user-defined data types, the local server must have compatible types of the same names.

Importing behavior can be customized with the following options (given in the `IMPORT FOREIGN SCHEMA` command):

`import_collate` (boolean)

This option controls whether column `COLLATE` options are included in the definitions of foreign tables imported from a foreign server. The default is `true`. You might need to turn this off if the remote

server has a different set of collation names than the local server does, which is likely to be the case if it's running on a different operating system. If you do so, however, there is a very severe risk that the imported table columns' collations will not match the underlying data, resulting in anomalous query behavior.

Even when this parameter is set to `true`, importing columns whose collation is the remote server's default can be risky. They will be imported with `COLLATE "default"`, which will select the local server's default collation, which could be different.

`import_default (boolean)`

This option controls whether column `DEFAULT` expressions are included in the definitions of foreign tables imported from a foreign server. The default is `false`. If you enable this option, be wary of defaults that might get computed differently on the local server than they would be on the remote server; `nextval()` is a common source of problems. The `IMPORT` will fail altogether if an imported default expression uses a function or operator that does not exist locally.

`import_generated (boolean)`

This option controls whether column `GENERATED` expressions are included in the definitions of foreign tables imported from a foreign server. The default is `true`. The `IMPORT` will fail altogether if an imported generated expression uses a function or operator that does not exist locally.

`import_not_null (boolean)`

This option controls whether column `NOT NULL` constraints are included in the definitions of foreign tables imported from a foreign server. The default is `true`.

Note that constraints other than `NOT NULL` will never be imported from the remote tables. Although Postgres Pro does support check constraints on foreign tables, there is no provision for importing them automatically, because of the risk that a constraint expression could evaluate differently on the local and remote servers. Any such inconsistency in the behavior of a check constraint could lead to hard-to-detect errors in query optimization. So if you wish to import check constraints, you must do so manually, and you should verify the semantics of each one carefully.

Tables or foreign tables which are partitions of some other table are imported only when they are explicitly specified in `LIMIT TO` clause.

F.11.1.10. Connection Management Options

By default, all connections that `postgres_fdw` establishes to foreign servers are kept open in the local session for re-use.

`keep_connections (boolean)`

This option controls whether `postgres_fdw` keeps the connections to the foreign server open so that subsequent queries can re-use them. It can only be specified for a foreign server. The default is `on`. If set to `off`, all connections to this foreign server will be discarded at the end of each transaction.

F.11.2. Functions

`postgres_fdw_get_connections(OUT server_name text, OUT valid boolean) returns setof record`

This function returns the foreign server names of all the open connections that `postgres_fdw` established from the local session to the foreign servers. It also returns whether each connection is valid or not. `false` is returned if the foreign server connection is used in the current local transaction but its foreign server or user mapping is changed or dropped (Note that server name of an invalid connection will be `NULL` if the server is dropped), and then such invalid connection will be closed at the end of that transaction. `true` is returned otherwise. If there are no open connections, no record is returned. Example usage of the function:

```
postgres=# SELECT * FROM postgres_fdw_get_connections() ORDER BY 1;
```

server_name	valid
-----+-----	
loopback1	t
loopback2	f

`postgres_fdw_disconnect(server_name text)` returns boolean

This function discards the open connections that are established by `postgres_fdw` from the local session to the foreign server with the given name. Note that there can be multiple connections to the given server using different user mappings. If the connections are used in the current local transaction, they are not disconnected and warning messages are reported. This function returns `true` if it disconnects at least one connection, otherwise `false`. If no foreign server with the given name is found, an error is reported. Example usage of the function:

```
postgres=# SELECT postgres_fdw_disconnect('loopback1');
postgres_fdw_disconnect
-----
t
```

`postgres_fdw_disconnect_all()` returns boolean

This function discards all the open connections that are established by `postgres_fdw` from the local session to foreign servers. If the connections are used in the current local transaction, they are not disconnected and warning messages are reported. This function returns `true` if it disconnects at least one connection, otherwise `false`. Example usage of the function:

```
postgres=# SELECT postgres_fdw_disconnect_all();
postgres_fdw_disconnect_all
-----
t
```

F.11.3. Connection Management

`postgres_fdw` establishes a connection to a foreign server during the first query that uses a foreign table associated with the foreign server. By default this connection is kept and re-used for subsequent queries in the same session. This behavior can be controlled using `keep_connections` option for a foreign server. If multiple user identities (user mappings) are used to access the foreign server, a connection is established for each user mapping.

When changing the definition of or removing a foreign server or a user mapping, the associated connections are closed. But note that if any connections are in use in the current local transaction, they are kept until the end of the transaction. Closed connections will be re-established when they are necessary by future queries using a foreign table.

Once a connection to a foreign server has been established, it's by default kept until the local or corresponding remote session exits. To disconnect a connection explicitly, `keep_connections` option for a foreign server may be disabled, or `postgres_fdw_disconnect` and `postgres_fdw_disconnect_all` functions may be used. For example, these are useful to close connections that are no longer necessary, thereby releasing connections on the foreign server.

F.11.4. Transaction Management

During a query that references any remote tables on a foreign server, `postgres_fdw` opens a transaction on the remote server if one is not already open corresponding to the current local transaction. The remote transaction is committed or aborted when the local transaction commits or aborts. Savepoints are similarly managed by creating corresponding remote savepoints.

The remote transaction uses `SERIALIZABLE` isolation level when the local transaction has `SERIALIZABLE` isolation level; otherwise it uses `REPEATABLE READ` isolation level. This choice ensures that if a query performs multiple table scans on the remote server, it will get snapshot-consistent results for all the scans. A consequence is that successive queries within a single transaction will see the same data from

the remote server, even if concurrent updates are occurring on the remote server due to other activities. That behavior would be expected anyway if the local transaction uses `SERIALIZABLE` or `REPEATABLE READ` isolation level, but it might be surprising for a `READ COMMITTED` local transaction. A future Postgres Pro release might modify these rules.

Note that it is currently not supported by `postgres_fdw` to prepare the remote transaction for two-phase commit.

F.11.5. Remote Query Optimization

`postgres_fdw` attempts to optimize remote queries to reduce the amount of data transferred from foreign servers. This is done by sending query `WHERE` clauses to the remote server for execution, and by not retrieving table columns that are not needed for the current query. To reduce the risk of misexecution of queries, `WHERE` clauses are not sent to the remote server unless they use only data types, operators, and functions that are built-in or belong to an extension that's listed in the foreign server's `extensions` option. Operators and functions in such clauses must be `IMMUTABLE` as well. For an `UPDATE` or `DELETE` query, `postgres_fdw` attempts to optimize the query execution by sending the whole query to the remote server if there are no query `WHERE` clauses that cannot be sent to the remote server, no local joins for the query, no row-level local `BEFORE` or `AFTER` triggers or stored generated columns on the target table, and no `CHECK OPTION` constraints from parent views. In `UPDATE`, expressions to assign to target columns must use only built-in data types, `IMMUTABLE` operators, or `IMMUTABLE` functions, to reduce the risk of misexecution of the query.

When `postgres_fdw` encounters a join between foreign tables on the same foreign server, it sends the entire join to the foreign server, unless for some reason it believes that it will be more efficient to fetch rows from each table individually, or unless the table references involved are subject to different user mappings. While sending the `JOIN` clauses, it takes the same precautions as mentioned above for the `WHERE` clauses.

The query that is actually sent to the remote server for execution can be examined using `EXPLAIN VERBOSE`.

F.11.6. Remote Query Execution Environment

In the remote sessions opened by `postgres_fdw`, the `search_path` parameter is set to just `pg_catalog`, so that only built-in objects are visible without schema qualification. This is not an issue for queries generated by `postgres_fdw` itself, because it always supplies such qualification. However, this can pose a hazard for functions that are executed on the remote server via triggers or rules on remote tables. For example, if a remote table is actually a view, any functions used in that view will be executed with the restricted search path. It is recommended to schema-qualify all names in such functions, or else attach `SET search_path` options (see [CREATE FUNCTION](#)) to such functions to establish their expected search path environment.

`postgres_fdw` likewise establishes remote session settings for various parameters:

- `TimeZone` is set to `UTC`
- `DateStyle` is set to `ISO`
- `IntervalStyle` is set to `postgres`
- `extra_float_digits` is set to 3 for remote servers 9.0 and newer and is set to 2 for older versions

These are less likely to be problematic than `search_path`, but can be handled with function `SET` options if the need arises.

It is *not* recommended that you override this behavior by changing the session-level settings of these parameters; that is likely to cause `postgres_fdw` to malfunction.

F.11.7. Cross-Version Compatibility

`postgres_fdw` can be used with remote servers dating back to PostgreSQL 8.3. Read-only capability is available back to 8.1. A limitation however is that `postgres_fdw` generally assumes that immutable

built-in functions and operators are safe to send to the remote server for execution, if they appear in a `WHERE` clause for a foreign table. Thus, a built-in function that was added since the remote server's release might be sent to it for execution, resulting in “function does not exist” or a similar error. This type of failure can be worked around by rewriting the query, for example by embedding the foreign table reference in a sub-`SELECT` with `OFFSET 0` as an optimization fence, and placing the problematic function or operator outside the sub-`SELECT`.

F.11.8. Wait Events

`postgres_fdw` can report the following wait events under the wait event type `Extension`:

`PostgresFdwCleanupResult`

Waiting for transaction abort on remote server.

`PostgresFdwConnect`

Waiting to establish a connection to a remote server.

`PostgresFdwGetResult`

Waiting to receive the results of a query from a remote server.

F.11.9. Configuration Parameters

`postgres_fdw.application_name` (string)

Specifies a value for [application_name](#) configuration parameter used when `postgres_fdw` establishes a connection to a foreign server. This overrides `application_name` option of the server object. Note that change of this parameter doesn't affect any existing connections until they are re-established.

`postgres_fdw.application_name` can be any string of any length and contain even non-ASCII characters. However when it's passed to and used as `application_name` in a foreign server, note that it will be truncated to less than `NAMEDATALEN` characters. Anything other than printable ASCII characters are replaced with [C-style hexadecimal escapes](#). See [application_name](#) for details.

`%` characters begin “escape sequences” that are replaced with status information as outlined below. Unrecognized escapes are ignored. Other characters are copied straight to the application name. Note that it's not allowed to specify a plus/minus sign or a numeric literal after the `%` and before the option, for alignment and padding.

Escape	Effect
<code>%a</code>	Application name on local server
<code>%c</code>	Session ID on local server (see log_line_prefix for details)
<code>%C</code>	Cluster name on local server (see cluster_name for details)
<code>%u</code>	User name on local server
<code>%d</code>	Database name on local server
<code>%p</code>	Process ID of backend on local server
<code>%%</code>	Literal <code>%</code>

For example, suppose user `local_user` establishes a connection from database `local_db` to foreign_db as user `foreign_user`, the setting `'db=%d, user=%u'` is replaced with `'db=local_db, user=local_user'`.

F.11.10. Examples

Here is an example of creating a foreign table with `postgres_fdw`. First install the extension:


```
CREATE EXTENSION postgres_fdw;
```

Then create a foreign server using [CREATE SERVER](#). In this example we wish to connect to a Postgres Pro server on host 192.83.123.89 listening on port 5432. The database to which the connection is made is named `foreign_db` on the remote server:

```
CREATE SERVER foreign_server
    FOREIGN DATA WRAPPER postgres_fdw
    OPTIONS (host '192.83.123.89', port '5432', dbname 'foreign_db');
```

A user mapping, defined with `CREATE USER MAPPING`, is needed as well to identify the role that will be used on the remote server:

```
CREATE USER MAPPING FOR local_user
    SERVER foreign_server
    OPTIONS (user 'foreign_user', password 'password');
```

Now it is possible to create a foreign table with `CREATE FOREIGN TABLE`. In this example we wish to access the table named `some_schema.some_table` on the remote server. The local name for it will be `foreign_table`:

```
CREATE FOREIGN TABLE foreign_table (
    id integer NOT NULL,
    data text
)
    SERVER foreign_server
    OPTIONS (schema_name 'some_schema', table_name 'some_table');
```

It's essential that the data types and other properties of the columns declared in `CREATE FOREIGN TABLE` match the actual remote table. Column names must match as well, unless you attach `column_name` options to the individual columns to show how they are named in the remote table. In many cases, use of `IMPORT FOREIGN SCHEMA` is preferable to constructing foreign table definitions manually.

F.11.11. Author

Shigeru Hanada <shigeru.hanada@gmail.com>

F.12. ptrack — a block-level incremental backup engine for Postgres Pro

PTRACK is a block-level incremental backup engine for Postgres Pro. If PTRACK is enabled, backup tools like [pg_probackup](#) can use its API to get information on the changed blocks on the fly when taking incremental backups. Copying only those blocks that have changed since the previous backup can significantly speed up the creation and minimize the size of backups.

PTRACK is designed to allow false positives, but to never allow false negatives: it means that all changes within the data directory except hint bits are guaranteed to be marked in the PTRACK map file, although some unchanged blocks might be included as well.

F.12.1. Setting up PTRACK

Once you complete Postgres Pro Shardman installation, do the following:

1. Add `ptrack` to the [shared_preload_libraries](#) parameter in the `postgresql.conf` file:

```
shared_preload_libraries = 'ptrack'
```

2. Set [ptrack.map_size](#) parameter to a positive integer.

For optimal performance, it is recommended to set `ptrack.map_size` to $N / 1024$, where N is the size of the Postgres Pro cluster, in MB. If you set this parameter to a lower value, PTRACK is more likely to map several blocks together, which leads to false-positive results when tracking changed blocks and increases the incremental backup size as unchanged blocks can also be copied into the incremental backup.

Setting `ptrack.map_size` to a higher value does not affect PTRACK operation, but keep in mind that you need up to `ptrack.map_size * 3` MB of additional disk space since PTRACK uses two additional temporary files to ensure durability. It is not recommended to set this parameter to a value higher than 1024. Even for Postgres Pro clusters with $N > 1$ TB, it makes sense to use a higher value only if the efficiency and sizes of backups are more critical than the potential overhead of maintaining a large PTRACK map.

3. Check the [wal_level](#) setting. When using PTRACK, it is required to set the [wal_level](#) parameter to `replica` or higher. Otherwise, you can lose some tracked changes if crash-recovery occurs: some commands do not write WAL at all if `wal_level` is `minimal`, and PTRACK map files are flushed to disk only at checkpoint time.
4. Restart the Postgres Pro Shardman instance for the changes to take effect, and then create the PTRACK extension:

```
CREATE EXTENSION ptrack;
```

As a result, several PTRACK functions are created, which are required for accessing PTRACK data.

Once this setup is complete, PTRACK starts tracking all the page changes in the Postgres Pro cluster and creates a `ptrack.map` file that stores the latest LSN values for these pages.

Note

The `ptrack.map_size` parameter can only be set at server start. If you change this parameter, the previously created PTRACK map file is cleared, and tracking newly changed blocks starts from scratch. To avoid losing recent changes, it is recommended to retake a full backup after modifying this setting.

F.12.2. PTRACK Configuration Parameters

`ptrack.map_size (integer)`

Specifies the size of a PTRACK map file and the amount of shared memory allocated for this file, in MB. The PTRACK map file stores the latest LSN values for all pages of the Postgres Pro cluster that have changed since PTRACK was enabled. It is not recommended to set this parameter to a value higher than 1GB. The 0 value disables PTRACK and cleans up all the related service files.

The `ptrack.map_size` parameter can only be set at server start. If you change this parameter, the previously created PTRACK map file is cleared, and tracking newly changed blocks starts from scratch. To avoid losing recent changes, it is recommended to retake a full backup after modifying this setting.

Default: 0

F.12.3. PTRACK Functions

`ptrack_init_lsn()` returns `pg_lsn`

Returns the LSN of the last PTRACK map initialization.

`ptrack_get_pagemapset (start_lsn pg_lsn)` returns `setof record`

Returns a list of data files changed since the specified `start_lsn` with the number and bitmap of changed pages for each file.

For example:

```
postgres=# SELECT * FROM ptrack_get_pagemapset ('0/185C8C0');
      path      | pagecount |          pagemap
-----+-----+-----
base/16384/1255 |          3 | \x00100000000050000000000000
base/16384/2674 |          3 | \x0000000900010000000000000000
base/16384/2691 |          1 | \x0000400000000000000000000000
base/16384/2608 |          1 | \x0000000000000000400000000000000000
base/16384/2690 |          1 | \x0004000000000000000000000000
```

`ptrack_get_change_stat (start_lsn pg_lsn)` returns `record`

Returns the statistics of changes (number of files, number of changed pages and their total size in MB) since the specified `start_lsn`.

For example:

```
postgres=# SELECT * FROM ptrack_get_change_stat ('0/285C8C8');
 files | pages |      size, MB
-----+-----+-----
    20 |    25 | 0.19531250000000000000
```

`ptrack_version()` returns `text`

Returns PTRACK version.

F.13. uuid-ossdp — a UUID generator

The `uuid-ossdp` module provides functions to generate universally unique identifiers (UUIDs) using one of several standard algorithms. There are also functions to produce certain special UUID constants. This module is only necessary for special requirements beyond what is available in core Postgres Pro. See [Section 9.14](#) for built-in ways to generate UUIDs.

This module is considered “trusted”, that is, it can be installed by non-superusers who have `CREATE` privilege on the current database.

F.13.1. uuid-ossdp Functions

[Table F.28](#) shows the functions available to generate UUIDs. The relevant standards ITU-T Rec. X.667, ISO/IEC 9834-8:2005, and [RFC 4122](#) specify four algorithms for generating UUIDs, identified by the version numbers 1, 3, 4, and 5. (There is no version 2 algorithm.) Each of these algorithms could be suitable for a different set of applications.

Table F.28. Functions for UUID Generation

Function	Description
<code>uuid_generate_v1 () → uuid</code>	Generates a version 1 UUID. This involves the MAC address of the computer and a time stamp. Note that UUIDs of this kind reveal the identity of the computer that created the identifier and the time at which it did so, which might make it unsuitable for certain security-sensitive applications.
<code>uuid_generate_v1mc () → uuid</code>	Generates a version 1 UUID, but uses a random multicast MAC address instead of the real MAC address of the computer.
<code>uuid_generate_v3 (namespace uuid, name text) → uuid</code>	Generates a version 3 UUID in the given namespace using the specified input name. The namespace should be one of the special constants produced by the <code>uuid_ns_*</code> () functions shown in Table F.29 . (It could be any UUID in theory.) The name is an identifier in the selected namespace. For example: <pre>SELECT uuid_generate_v3(uuid_ns_url(), 'http://www.postgresql.org');</pre> The name parameter will be MD5-hashed, so the cleartext cannot be derived from the generated UUID. The generation of UUIDs by this method has no random or environment-dependent element and is therefore reproducible.
<code>uuid_generate_v4 () → uuid</code>	Generates a version 4 UUID, which is derived entirely from random numbers.
<code>uuid_generate_v5 (namespace uuid, name text) → uuid</code>	Generates a version 5 UUID, which works like a version 3 UUID except that SHA-1 is used as a hashing method. Version 5 should be preferred over version 3 because SHA-1 is thought to be more secure than MD5.

Table F.29. Functions Returning UUID Constants

Function	Description
<code>uuid_nil () → uuid</code>	Returns a “nil” UUID constant, which does not occur as a real UUID.
<code>uuid_ns_dns () → uuid</code>	Returns a constant designating the DNS namespace for UUIDs.

Function	Description
<code>uuid_ns_url</code> <code>() → uuid</code>	Returns a constant designating the URL namespace for UUIDs.
<code>uuid_ns_oid</code> <code>() → uuid</code>	Returns a constant designating the ISO object identifier (OID) namespace for UUIDs. (This pertains to ASN.1 OIDs, which are unrelated to the OIDs used in Postgres Pro.)
<code>uuid_ns_x500</code> <code>() → uuid</code>	Returns a constant designating the X.500 distinguished name (DN) namespace for UUIDs.

F.13.2. Building `uuid-oss`

Historically this module depended on the OSSP UUID library, which accounts for the module's name. While the OSSP UUID library can still be found at <http://www.ossf.org/pkg/lib/uuid/>, it is not well maintained, and is becoming increasingly difficult to port to newer platforms. `uuid-oss` can now be built without the OSSP library on some platforms. On FreeBSD and some other BSD-derived platforms, suitable UUID creation functions are included in the core `libc` library. On Linux, macOS, and some other platforms, suitable functions are provided in the `libuuid` library, which originally came from the `e2fsprogs` project (though on modern Linux it is considered part of `util-linux-ng`). When invoking `configure`, specify `--with-uuid=bsd` to use the BSD functions, or `--with-uuid=e2fs` to use `e2fsprogs`' `libuuid`, or `--with-uuid=oss` to use the OSSP UUID library. More than one of these libraries might be available on a particular machine, so `configure` does not automatically choose one.

F.13.3. Author

Peter Eisentraut <peter_e@gmx.net>

Appendix G. Postgres Pro Modules and Extensions Shipped as Individual Packages

This appendix contains Postgres Pro modules and extensions that are made available in Postgres Pro Shardman as individual packages. These packages are listed in [Section 17.1.4.1](#), and the documentation for each module explains which package to choose to install it. [Appendix F](#) covers more modules and extensions included in the Postgres Pro Shardman distribution. Note also the `pg_portal_modify` extension, which is shipped as an individual package, and the `pgpro-orautl` package, which provides packages for `utl_http`, `utl_mail`, and `utl_smtp` (see [Section 17.1.4.1](#) for details).

G.1. pgpro_pwr — workload reports

The `pgpro_pwr` module is designed to discover most resource-intensive activities in your database. (PWR, pronounced like "power", is an abbreviation of Postgres Pro Workload Reporting.) This extension is based on Postgres Pro's [Statistics Collector](#) views and the `pgpro_stats` or `pg_stat_statements` extension.

Note

Although `pgpro_pwr` can work with the `pg_stat_statements` extension, it is recommended that you use the `pgpro_stats` extension since it provides statement plans, wait events sampling and load distribution statistics for databases, roles, client hosts and applications.

Below, use of `pgpro_stats` is assumed unless otherwise noted.

If you cannot use `pgpro_stats` for an observed database, but the `pg_stat_kcache` extension is available, `pgpro_pwr` can process `pg_stat_kcache` data, which also provides information about CPU resource usage of statements and filesystem load (*rusage*).

`pgpro_pwr` is based on cumulative statistics sampling. Each sample contains statistic increments for most active objects and queries since the time when the previous sample was taken, or more concisely, *since the previous sample*. This data is later used to generate reports.

`pgpro_pwr` provides functions to collect samples. Regular sampling allows building a report on the database workload in the past.

`pgpro_pwr` allows you to take explicit samples during batch processing, load testing, etc.

Any time a sample is taken, `pgpro_stats_statements_reset()` (see [pgpro_stats](#) for the function description) is called to ensure that statement statistics will not be lost when the statements count exceed `pgpro_stats.max` (see [Section F.7.7.1](#)). The report will also contain a section informing you of whether the count of captured statements in any sample reaches 90% of `pgpro_stats.max`.

`pgpro_pwr` installed on one Postgres Pro server can also collect statistics from other servers. This feature is useful for gathering workload statistics from hot standbys on the primary server. To benefit from it, make sure that all server names and connection strings are specified and that the `pgpro_pwr` server can connect to all databases on all servers.

G.1.1. pgpro_pwr Architecture

The extension consists of the following parts:

- *Historical repository* is a storage for sampling data. The repository is a set of extension tables.

Note

Among the rest, `pgpro_pwr` tables store query texts, which can contain sensitive information. So, for security reasons, restrict access to the repository as appropriate.

- *Sample management engine* comprises functions used to [take samples](#) and maintain the repository by removing obsolete sample data.
- *Report engine* comprises functions for [generating reports](#) based on data from the historical repository.
- *Administrative functions* allow you to create and manage [servers](#) and [baselines](#).

G.1.2. Prerequisites

The prerequisites assume that `pgpro_pwr`, which is usually installed in a *target* cluster, i.e., the cluster that you will mainly track the workload for, the extension can also collect performance data from other clusters.

G.1.2.1. For the `pgpro_pwr` Database

The `pgpro_pwr` extension depends on [PL/pgSQL](#) and the [dblink](#) extension.

G.1.2.2. For the Target Server

The target server must allow connections to all databases from the server where `pgpro_pwr` is running. To connect to the target server, provide a connection string where a particular database on this server is specified. This database is of high importance for `pgpro_pwr` since the functionality of the [pgpro_stats](#) or [pg_stat_statements](#) extensions will be provided through this database. Note, however, that `pgpro_pwr` will also connect to all the other databases on this server.

Optionally, for completeness of gathered statistics:

- If statement statistics are needed in reports, [pgpro_stats](#) must be installed and configured in the aforementioned database. The following settings may affect the completeness and accuracy of gathered statistics:

- `pgpro_stats.max`

Low setting of this parameter may cause some statement statistics to be wiped out before the sample is taken. A report will warn you if the value of `pgpro_stats.max` seems undersized.

- `pgpro_stats.track`

Avoid changing the default value of 'top' (note that the value of 'all' will affect the accuracy of %Total fields for statements-related sections of a report).

- Set the parameters of the Postgres Pro's [Statistics Collector](#) as follows:

```
track_activities = on
track_counts = on
track_io_timing = on
track_wal_io_timing = on    # Since PostgreSQL 14
track_functions = all/pl
```

G.1.3. Installation and Setup

`pgpro_pwr` is provided with Postgres Pro Shardman as a separate pre-built package `pgpro-pwr-sdm-17` (for the detailed installation instructions, see [Chapter 17](#)).

Note

`pgpro_pwr` creates a bunch of database objects, so installation in a dedicated schema is recommended.

Although the use of `pgpro_pwr` with superuser privileges does not have any issues, superuser privileges are not necessary. So you can choose one of the following setup procedures depending on your configuration and security requirements or customize them to meet your needs:

G.1.3.1. Simple Setup

Use this setup procedure when `pgpro_pwr` is to be installed on the target cluster to only track its workload as superuser.

Create a schema for the `pgpro_pwr` installation and create the extension:

```
CREATE SCHEMA profile;
CREATE EXTENSION pgpro_pwr SCHEMA profile;
```

G.1.3.2. Complex Setup

Use this setup procedure when you intend to use `pgpro_pwr` for tracking workload on one or more servers and need to follow the principle of least privilege.

G.1.3.2.1. In the Target Server Database

Create a user for `pgpro_pwr` on the target server:

```
CREATE USER pwr_collector PASSWORD 'collector_pwd';
```

Make sure this user has permissions to connect to any database in the target cluster (by default, it is true) and that `pg_hba.conf` permits such a connection from the `pgpro_pwr` database host. Also, grant `pwr_collector` with membership in the `pg_read_all_stats` role and the `EXECUTE` privilege on the following functions:

```
GRANT pg_read_all_stats TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgpro_stats_statements_reset TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgpro_stats_totals_reset(text,bigint) TO pwr_collector;
```

Also ensure the `SELECT` privilege on the `pgpro_stats_archiver` view:

```
GRANT SELECT ON pgpro_stats_archiver TO pwr_collector;
```

G.1.3.2.2. In the `pgpro_pwr` Database

Create an unprivileged user:

```
CREATE USER pwr_user;
```

This user will be the owner of the extension schema and will collect samples.

Create a schema for the `pgpro_pwr` installation:

```
CREATE SCHEMA profile AUTHORIZATION pwr_user;
```

Grant the `USAGE` privilege on the schema where the `dblink` extension resides:

```
GRANT USAGE ON SCHEMA public TO pwr_user;
```

Create the extension using `pwr_user` account:

```
\c - pwr_user
```



```
CREATE EXTENSION pgpro_pwr SCHEMA profile;
```

Define the connection parameters of the target server for pgpro_pwr. For example:

```
SELECT profile.create_server('target_server_name', 'host=192.168.1.100 dbname=postgres  
port=5432');
```

The connection string provided will be used in the `dblink_connect()` call while executing the `take_sample()` function.

Note

Connection strings are stored in a pgpro_pwr table in clear-text form. Make sure no other database users can access tables of the pgpro_pwr extension.

G.1.3.3. Setting Up pgpro_pwr Roles

Up to three roles can be distinguished when pgpro_pwr is in operation:

- The *pgpro_pwr owner* role is the owner of the pgpro_pwr extension.
- The *collecting role* is used by pgpro_pwr to connect to databases and collect statistics.
- The *reporting role* is used to generate reports.

If all the actions with pgpro_pwr are performed by the superuser role `postgres`, you can skip most of the setup explained below.

G.1.3.3.1. The pgpro_pwr Owner

This role can be used to perform all actions related to pgpro_pwr. This role will have access to server connection strings, which may contain passwords. You should use this role to call the `take_sample()` function. The `dblink` extension is needed for this user.

Consider an example assuming each extension in its own schema:

```
\c postgres postgres
CREATE SCHEMA dblink;
CREATE EXTENSION dblink SCHEMA dblink;
CREATE USER pwr_usr with password 'pwr_pwd';
GRANT USAGE ON SCHEMA dblink TO pwr_usr;
CREATE SCHEMA profile AUTHORIZATION pwr_usr;
\c postgres pwr_usr
CREATE EXTENSION pgpro_pwr SCHEMA profile;
```

G.1.3.3.2. The Collecting Role

This role should be used by pgpro_pwr to connect to databases and collect statistics. Unprivileged users cannot open connections using dblink without a password, so you need to provide the password in the connection string for each server. This role should have access to all supported statistics extensions. It should also be able to perform a reset of statistics extensions.

Consider an example. If you use `pgpro_stats` to collect statistics, set up the collecting role as follows:

```
\c postgres postgres
CREATE SCHEMA pgps;
CREATE EXTENSION pgpro_stats SCHEMA pgps;
CREATE USER pwr_collector with password 'collector_pwd';
GRANT pg_read_all_stats TO pwr_collector;
GRANT USAGE ON SCHEMA pgps TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgps.pgpro_stats_statements_reset TO pwr_collector;
```

If you use `pg_stat_statements` to collect statistics, set up the collecting role as follows:

```
\c postgres postgres
CREATE SCHEMA pgss;
CREATE SCHEMA pgsk;
CREATE SCHEMA pgws;
CREATE EXTENSION pg_stat_statements SCHEMA pgss;
CREATE EXTENSION pg_stat_kcache SCHEMA pgsk;
CREATE EXTENSION pg_wait_sampling SCHEMA pgws;
CREATE USER pwr_collector with password 'collector_pwd';
GRANT pg_read_all_stats TO pwr_collector;
GRANT USAGE ON SCHEMA pgss TO pwr_collector;
GRANT USAGE ON SCHEMA pgsk TO pwr_collector;
GRANT USAGE ON SCHEMA pgws TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgss.pg_stat_statements_reset TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgsk.pg_stat_kcache_reset TO pwr_collector;
GRANT EXECUTE ON FUNCTION pgws.pg_wait_sampling_reset_profile TO pwr_collector;
```

Now you should set up a connection string pointing to the database with statistics extensions installed:

```
\c postgres pwr_usr
SELECT profile.set_server_connstr('local','dbname=postgres port=5432 host=localhost
user=pwr_collector password=collector_pwd');
```

Password authentication must be configured in the `pg_hba.conf` file for the `pwr_collector` user.

Obviously, the collecting role should be properly configured on all servers observed by `pgpro_pwr`.

Now you should be able to call `take_sample()` using the `pwr_usr` role:

```
\c postgres pwr_usr
SELECT * FROM take_sample();
```

And now it's time to configure the scheduler (in our example, the `crontab` command of the `postgres` user):

```
*/30 * * * *    psql -U pwr_usr -d postgres -c 'SELECT profile.take_sample()' > /dev/
null 2>&1
```

Note that you can use the Postgres Pro password file to store passwords.

G.1.3.3.3. The Reporting Role

Any user can build a `pgpro_pwr` report. The minimal privileges needed to generate `pgpro_pwr` reports are granted to the `public` role. However a full report, with query texts, is only available to the member of the `pg_read_all_stats` role. Anyway, the reporting role cannot access server connection strings, so it cannot get the passwords of servers.

G.1.3.3.4. Setting Extension Parameters

In `postgresql.conf`, you can define the following `pgpro_pwr` parameters:

`pgpro_pwr.max (integer)`

Number of top objects (statements, relations, etc.) to be reported in each sorted report table. This parameter affects the size of a sample.

The default value is 20.

`pgpro_pwr.max_sample_age (integer)`

Retention time of the sample, in days. Samples aged `pgpro_pwr.max_sample_age` days and older are automatically deleted on the next `take_sample()` call.

The default value is 7 days.

`pgpro_pwr.max_query_length` (integer)

Maximum query length allowed in reports. All queries in a report will be truncated to `pgpro_pwr.max_query_length` characters.

The default value is 20 000 characters.

`pgpro_pwr.track_sample_timings` (boolean)

Enables collecting detailed timing statistics of `pgpro_pwr`'s own sampling procedures. Set this parameter to diagnose why sampling functions run slowly. Collected timing statistics will be available in the `v_sample_timings` view.

The default value is `off`.

`pgpro_pwr.statements_reset` (boolean)

Controls the `pgpro_stats/pg_stat_statements` statistics reset during taking a sample. Allows you not to reset the statistics during taking a sample due to new techniques employed. When disabled, `pgpro_pwr` will track statement evictions using the value of the `calls` field. However this method does not completely prevent statistics loss. `pg_stat_statements` v1.11 and `pgpro_stats` v1.8 contain time tracking abilities that can reduce the possible data loss. When this setting is disabled, you can temporarily enable it in a session if you want to sometimes perform a reset of `pgpro_stats/pg_stat_statements`.

The default value is `on`.

`pgpro_pwr.relsizes_collect_mode` (text)

Defines the mode of collecting relation sizes. Possible values:

- `off` — collection of relation sizes is based on `pg_class`. Although the relation sizes collected this way are rough, their collection consumes almost no resources.
- `on` — accurate relation sizes are collected for each sample using the `pg_relation_size()` function, which requires a lock on the table and is pretty resource intensive.
- `schedule` — accurate relation sizes are collected in the size-collection window, defined for each server.

The default value is `off`.

G.1.4. Managing Servers

Once installed, `pgpro_pwr` creates one enabled `local` server for the current cluster. If a server is *enabled*, `pgpro_pwr` includes it in sampling when no server is explicitly specified (see `take_sample()` for details). A server that is not enabled is referred to as *disabled*.

The default connection string for a local node contains only `dbname` and `port` parameters. The values of these parameters are taken from the connection used to create the extension. You can change the server connection string using the `set_server_connstr()` function when needed.

G.1.4.1. Server Management Functions

Use the following `pgpro_pwr` functions for server management:

```
create_server(server name, connstr text, enabled boolean DEFAULT TRUE, max_sample_age
integer DEFAULT NULL description text DEFAULT NULL)
```

Creates a server definition.

Arguments:

- `server` — server name. Must be unique.

- *connstr* — connection string. Must contain all the necessary settings to connect from pgpro_pwr server to the target server database.
- *enabled* — set to include the server in sampling by the `take_sample()` function without arguments.
- *max_sample_age* — retention time of the sample. Overrides the global `pgpro_pwr.max_sample_age` setting for this server.
- *description* — server description text, to be included in reports.

Here is an example of how to create a server definition:

```
SELECT profile.create_server('omega','host=192.168.1.100 dbname=postgres
port=5432');
```

`drop_server(server name)`

Drops a server and all its samples.

`set_server_description(server name description text)`

Sets a new server description.

`set_server_subsampling(server name, subsample_enabled boolean, min_query_duration interval, min_xact_duration interval, min_xact_age integer, min_idle_xact_dur interval hour to second)`

Defines subsample settings for a server.

Arguments:

- *server* — server name.
- *subsample_enabled* — defines whether subsampling is enabled for the server, that is, whether `take_subsample()` function should actually take a subsample.
- *min_query_duration* — the query duration threshold.
- *min_xact_duration* — the transaction duration threshold.
- *min_xact_age* — the transaction age threshold.
- *min_idle_xact_dur_age* — the idle transaction threshold.

`enable_server(server name)`

Includes a server in sampling by the `take_sample()` function without arguments.

`disable_server(server name)`

Excludes a server from sampling by the `take_sample()` function without arguments.

`rename_server(server name, new_name name)`

Renames a server.

`set_server_max_sample_age(server name, max_sample_age integer)`

Sets the retention period for a server (in days). To reset the server retention, set the value of *max_sample_age* to NULL.

`set_server_db_exclude(server name, exclude_db name[])`

Excludes a list of databases on a server from sampling. Use when pgpro_pwr is unable to connect to some databases in a cluster (for example, in Amazon RDS instances).

`set_server_connstr(server name, server_connstr text)`

Sets the connection string for a server.

`set_server_setting(server name, setting text, value jsonb)`

Fine-tunes settings of the server statistics collection. `collect*` settings control which statistics should be collected, and `value` for these settings accepts boolean values, with the default equal to `true`. Available settings:

- `collect_pg_stat_statement` — collect statement statistics using `pg_stat_statements` and `pg_stat_kcache` extensions.
- `collect_objects` — collect all the schema object statistics, that is, tables, indexes, and functions, from `pg_stat_*` views.
- `collect_relations` — collect statistics on tables and indexes from `pg_stat_*` views.
- `collect_functions` — collect statistics on user functions from the `pg_stat_user_functions` view.
- `collect_vacuum_stats` — collect the extended vacuum statistics.

`show_server_settings(server name)`

Returns the statistics collection settings for the specified server.

`show_servers()`

Displays the list of configured servers.

G.1.5. Managing Samples

A *sample* contains the database workload statistics since the previous sample

G.1.5.1. Sampling Functions

The following `pgpro_pwr` functions relate to sampling:

`take_sample()`

`take_sample(server name [, skip_sizes boolean])`

Takes samples.

If the parameter is omitted, the function takes a sample on each enabled server. Servers are accessed for sampling sequentially, one by one. The function returns a table with the following columns:

- `server` — server name.
- `result` — result of taking the sample. Can be `OK` if the sample was taken successfully or contain the error trace text in case of exception.
- `elapsed` — time elapsed while the sample was taken.

If called with the parameter, the function takes a sample on the specified server even if this server is disabled. Use when you need different sampling frequencies on specific servers. Returns 0 on success.

Arguments:

- `server` — server name.
- `skip_sizes` — if omitted or set to null, the size-collection policy applies; if `false`, relation sizes are collected; if `true`, the collection of relation sizes is skipped.

`take_sample_subset([sets_cnt integer, current_set integer])`

Takes a sample on each server in a subset of servers. Use to take samples on servers in parallel if you have many enabled servers. Although PL/pgSQL does not support parallel execution, you can call this function in parallel sessions. This function returns the same type as `take_sample()`. If both parameters are omitted, the function behaves like the `take_sample()` function, i.e., it takes a sample on all enabled servers one by one.

Arguments:

- *sets_cnt* — number of subsets to divide all enabled servers into.
- *current_set* — number of the subset to collect samples for. Takes values from 0 through *sets_cnt* - 1. For the specified subset, samples are collected as usual, server by server.

If a reset of statistics since the previous sample was detected, `pgpro_pwr` treats corresponding absolute values as differentials; however, the accuracy will be affected anyway.

```
show_samples([server name,] [days integer])
```

Returns a table with information on server samples (`local` server is assumed if *server* is omitted) for the last *days* days (all existing samples are assumed if omitted). This table has the following columns:

- *sample* — sample identifier.
- *sample_time* — time when this sample was taken.
- *dbstats_reset* — NULL or the statistics reset timestamp of the `pg_stat_database` view if the statistics were reset since the previous sample.
- *clustats_reset* — NULL or the statistics reset timestamp of the `pg_stat_bgwriter` view if the statistics were reset since the previous sample.
- *archstats_reset* — NULL or the statistics reset timestamp of the `pg_stat_archiver` view if the statistics were reset since the previous sample.

Sampling functions also maintain the server repository by deleting obsolete samples and baselines according to the retention policy.

G.1.5.2. Taking Samples

To take samples for all enabled servers, call the `take_sample()` function. Usually, one or two samples per hour is sufficient. You can use a cron-like tool to schedule sampling. Here is an example for a 30-minute sampling period:

```
*/30 * * * * psql -c 'SELECT profile.take_sample()' &> /dev/null
```

However, the results of such a call are not checked for errors. In a production environment, function results can be used for monitoring. This function returns OK for all servers with successfully taken samples and shows error text for failed servers:

```
SELECT * FROM take_sample();
  server  |                               result
-----+-----
| elapsed
+-----+-----
ok_node  | OK
| 00:00:00.48
fail_node | could not establish connection
+| 00:00:00
| SQL statement "SELECT dblink_connect('server_connection',server_connstr)"
+|
| PL/pgSQL function take_sample(integer) line 69 at PERFORM
+|
| PL/pgSQL function take_sample_subset(integer,integer) line 27 at
assignment+|
| SQL function "take_sample" statement 1
+|
| FATAL:  database "postgresno" does not exist
|
(2 rows)
```

G.1.5.3. Sample Retention Policy

You can define sample retention at the following levels:

1. *Global*

The value of the `pgpro_pwr.max_sample_age` parameter in the `postgresql.conf` file defines a common retention setting, which is effective if none of other related settings are defined.

2. *Server*

Specifying the `max_sample_age` parameter while creating a server or calling the `set_server_max_sample_age(server, max_sample_age)` function for an existing server defines the retention for the server. A server retention setting overrides `pgpro_pwr.max_sample_age` for a specific server.

3. *Baseline*

A [baseline](#) created overrides all the other retention periods for included samples.

G.1.6. Managing the Collection of Relation Sizes

It may take considerable time to collect sizes of all relations in a database by Postgres Pro relation-size functions. Besides, those functions require `AccessExclusiveLock` on a relation. However, it may be sufficient for you to collect relation sizes on a daily basis. `pgpro_pwr` allows you to skip collecting relation sizes by defining the *size-collection policy* for servers. The policy defines:

- A daily window when the collection of relation sizes is permitted.
- A minimum gap between two samples with relation sizes collected.

When the size-collection policy is defined, sampling functions collect relation sizes only when the sample is taken in the defined window and the previous sample with sizes is older than the gap. The following function defines this policy:

```
set_server_size_sampling(server name, window_start time with time zone DEFAULT NULL,  
window_duration interval hour to second DEFAULT NULL, sample_interval interval day to  
minute DEFAULT NULL, collect_mode text DEFAULT NULL)
```

Defines the size-collection policy for a server.

Arguments:

- `server` — server name.
- `window_start` — start time of the size-collection window.
- `window_duration` — duration of the size-collection window.
- `sample_interval` — minimum time gap between two samples with relation sizes collected.
- `collect_mode` — when set to `off`, which is the default for new installations, relation sizes are collected from `pg_class`, when set to `on`, relation sizes are collected using the `pg_relation_size()` function, when set to `schedule`, `pgpro_pwr` collects the relation size in the size-collection window. This parameter overrides the `relnsize_collect_mode` extension parameter. Upgrading from a previous version sets the value of this parameter to `on` or `schedule`, so the previous behavior does not change.

Note

When you build a report between samples either of which lacks relation-size data, relation-growth sections will be based on `pg_class.relpages` data. However, you can expand the report interval bounds to the nearest samples with relation sizes collected using the `with_growth` parameter of [report generation functions](#); this makes the growth data more accurate.

Relation sizes are needed to calculate sequentially scanned volume for tables and explicit vacuum load for indexes.

Example:

```
SELECT set_server_size_sampling('local','23:00+03',interval '2 hour',interval '8 hour',
'schedule');
```

The `show_servers_size_sampling` function shows size collection policies for all servers:

```
postgres=# SELECT * FROM show_servers_size_sampling();
 server_name | window_start | window_end | window_duration | sample_interval |
limited_collection
-----+-----+-----+-----+-----+-----
+-----+
 local      | 23:00:00+03 | 01:00:00+03 | 02:00:00        | 08:00:00        | t
```

G.1.7. Managing Subsamples

Some performance-related data available in Postgres Pro is not cumulative. For example, the most often used data about session states is available through the [pg_stat_activity](#) view and can only be obtained with frequent samples. However, the `take_sample()` function is heavy and can take considerable amount of time. So it is not suitable for collecting session state data.

The *subsample* feature provides a new fast `take_subsample()` function. It can be used to collect relatively fast changing data. Every subsample is bound to the next regular sample and is deleted by the retention policy together with it.

The subsample feature can be used to capture the most interesting session states:

- Long running queries
- Long transactions
- Aged transactions, that is, those that hold a snapshot behind a lot of other transactions
- Transactions being in an idle state for a long time

G.1.7.1. Subsampling Functions

The following pgpro_pwr functions relate to subsampling:

```
take_subsample()
take_subsample(server name)
```

If the parameter is omitted, the function takes a subsample on each *enabled* server with subsampling enabled (see [set_server_subsampling](#) for details). Server subsamples are taken sequentially, one by one. The function returns a table with the following columns:

- `server` — server name.
- `result` — result of taking the subsample. Can be `OK` if the subsample was taken successfully or contain the error text in case of exception.
- `elapsed` — time elapsed while the subsample was taken.

This tabular return format makes it easy to control subsample creation using an SQL query.

If called with the parameter, the function takes a subsample for the specified server. Use when you need different subsampling frequencies on servers or if you want to take an explicit subsample on a specific server.

Arguments:

- *server* — server name.

Note

Trying to take a subsample during taking a sample fails.

```
take_subsample_subset([sets_cnt integer], [current_set integer])
```

Takes subsamples for a subset of enabled servers with subsampling enabled. Although subsamples should be fast enough for serial processing, subsamples can be taken in parallel, like regular samples. This function returns the same type as `take_subsample()`. If both parameters are omitted, the function behaves like the `take_subsample()` function.

Arguments:

- *sets_cnt* — number of server subsets.
- *current_set* — the subset to process. Takes values from 0 through *sets_cnt* - 1. For the specified subset, subsamples are collected as usual, server by server.

G.1.7.2. Configuring the Subsample Feature

The following settings affect the subsample behaviour:

- `pgpro_pwr.subsample_enabled` — defines whether the `take_subsample()` function should actually take a subsample.
- `pgpro_pwr.min_query_duration` — the long running query threshold.
- `pgpro_pwr.min_xact_duration` — the long transaction threshold.
- `pgpro_pwr.min_xact_age` — the transaction age threshold.
- `pgpro_pwr.min_idle_xact_dur_age` — the idle transaction threshold.

The subsample behavior can be defined at the server level using the `set_server_subsampling` function.

The last observed session state is saved in the repository when either of the following threshold-related events happens:

- During query execution, the difference between `now()` and the `query_start` exceeds the `pgpro_pwr.min_query_duration` threshold.
- During a transaction, the difference between `now()` and the `xact_start` exceeds the `pgpro_pwr.min_xact_duration` threshold.
- During a transaction, the `age(backend_xmin)` exceeds the `pgpro_pwr.min_xact_age` threshold.
- During a transaction in a state `idle in transaction` or `idle in transaction (aborted)`, the difference between `now()` and the `state_change` exceeds the `pgpro_pwr.min_idle_xact_duration` threshold.

See [Chapter 27](#) for details of the mentioned fields. Every subsample can hold at most `pg_profile.topn` entries for every threshold type.

G.1.7.3. Scheduling Subsamples

Subsamples are fast enough to take them quite often. However usually you do not need more than 2-4 subsamples per minute. Obviously the subsample frequency depends on the shortest used duration of a threshold setting.

Cron only allows one call per minute, so some effort is needed to schedule more frequent subsamples. For example, the `\watch psql` command can be used:

```
echo "select take_subsample(); \watch 15" | psql &> /dev/null
```

The `psql` call can be wrapped in the `systemd` service like this:

```
Description=pg_profile subsampling unit
[Unit]

[Service]
Type=simple
ExecStart=/bin/sh -c 'echo "select take_subsample(); \\watch 15" | /path/to/psql -go /dev/null'
User=postgres
Group=postgres

[Install]
WantedBy=multi-user.target
```

G.1.8. Managing Baselines

A *baseline* is a named sequence of samples that has its own retention setting. A baseline can be used as a sample interval in report generation functions. An undefined baseline retention means infinite retention. Use baselines to save information about the database workload for a certain time interval.

G.1.8.1. Baseline Management Functions

Use the following `pgpro_pwr` functions for baseline management:

```
create_baseline([server name,] baseline varchar(25), start_id integer, end_id integer [,
days integer])
create_baseline([server name,] baseline varchar(25), time_range tstzrange [, days in-
teger])
```

Creates a baseline.

Arguments:

- *server* — server name. `local` sever is assumed if omitted.
- *baseline* — baseline name. Must be unique for a server.
- *start_id* — identifier of the first sample in the baseline.
- *end_id* — identifier of the last sample in the baseline.
- *time_range* — time interval for the baseline. The baseline will include all samples for the minimal interval that covers *time_range*.
- *days* — baseline retention time, defined in integer days since `now()`. Omit or set to null for infinite retention.

```
drop_baseline([server name,] baseline varchar(25))
```

Drops a baseline. For the meaning and usage details of function arguments, see [create_baseline](#). Dropping a baseline does not mean dropping all its samples immediately. The baseline retention just no longer applies to them.

```
keep_baseline([server name,] baseline varchar(25) [, days integer])
```

Changes the retention of a baseline. For the meaning and usage details of function arguments, see [create_baseline](#). Omit the *baseline* parameter or pass null to it to change the retention of all existing baselines.

```
show_baselines([server name])
```

Displays existing baselines. Call `show_baselines` to get information about the baselines, such as names, sampling intervals and retention periods. `local` sever is assumed if the *server* parameter is omitted.

G.1.9. Data Export and Import

Collected samples can be exported from one instance of the `pgpro_pwr` extension and then loaded into another one. This feature helps you to move server data from one instance to another or to send collected data to your support team.

G.1.9.1. Data Export

The `export_data` function exports data to a regular table. You can use any method available to export this table from your database. For example, you can use the `\copy` meta-command of `psql` to obtain a single `csv` file:

```
postgres=# \copy (select * from export_data()) to 'export.csv'
```

G.1.9.2. Data Import

Since data can only be imported from a local table, first, load the data you exported. Using the `\copy` meta-command again:

```
postgres=# CREATE TABLE import (section_id bigint, row_data json);
CREATE TABLE
postgres=# \copy import from 'export.csv'
COPY 6437
```

Now you can import the data by providing the `import` table to the `import_data` function:

```
postgres=# SELECT * FROM import_data('import');
```

After successful import, you can drop the `import` table.

Note

If server data is imported for the first time, your local `pgpro_pwr` servers with matching names will cause a conflict during import. To avoid this, you can temporarily rename such servers or you can specify the server name prefix for import operations. However, during import of new data for already imported servers, they are matched by system identifiers, so feel free to rename imported servers. Also keep in mind that `pgpro_pwr` sets servers being imported to the *disabled* state for `take_sample()` to bypass them.

G.1.9.3. Export and Import Functions

Use these functions to export or import data:

```
export_data([server name, [min_sample_id integer,] [max_sample_id integer,]] [, obfuscate_queries boolean])
```

Exports collected data.

Arguments:

- `server` — server name. All configured servers are assumed if omitted.
- `min_sample_id`, `max_sample_id` — sample identifiers to bound the export (inclusive). If `min_sample_id` is omitted or set to null, all samples until `max_sample_id` sample are exported; if `max_sample_id` is omitted or set to null, all samples since `min_sample_id` sample are exported.
- `obfuscate_queries` — if `true`, query texts are exported as MD5 hash.

```
import_data(data regclass [, server_name_prefix text])
```

Imports previously exported data. Returns the number of actually loaded rows in `pgpro_pwr` tables.

Arguments:

- *data* is the name of the table containing import data.
- *server_name_prefix* specifies the server name prefix for the import operation. It can be used to avoid name conflicts.

G.1.10. Report Generation Functions

pgpro_pwr reports are generated in HTML format by reporting functions. The following types of reports are available:

- *Regular reports* provide statistics on the workload for an interval.
- *Differential reports* provide statistics on the same objects for two intervals. Corresponding values are located next to each other, which makes it easy to compare the workloads.

Reporting functions take sample identifiers, baselines or time ranges to determine the intervals. For time ranges, these are the minimal intervals that cover the ranges.

G.1.10.1. Regular Reports

Use these functions to generate regular reports:

```
get_report([server name,] start_id integer, end_id integer [, description text [,
with_growth boolean [, db_exclude name[]]])
get_report([server name,] time_range tstzrange [, description text [, with_growth boolean
[, db_exclude name[]]])
get_report([server name,] baseline varchar(25) [, description text [, with_growth boolean
[, db_exclude name[]]])
```

Generates a regular report defined by the arguments.

Arguments:

- *server* — server name. `local` sever is assumed if omitted.
- *start_id* — identifier of the interval starting sample.
- *end_id* — identifier of the interval ending sample.
- *baseline* — baseline name.
- *time_range* — time range.
- *description* — short text to be included in the report as its description.
- *with_growth* — flag requesting interval expansion to the nearest bounds with data on relation growth available. The default value is `false`.
- *db_exclude* — the database exclusion list. Lists databases to be excluded from all report tables having the `Database` column. Use to hide some databases in the report.

```
get_report_latest([server name,])
get_report_latest([server name [, db_exclude name[]]])
```

Generates a regular report for two latest samples.

Arguments:

- *server* — server name. `local` sever is assumed if omitted.
- *db_exclude* — the database exclusion list. Lists databases to be excluded from all report tables having the `Database` column. Use to hide some databases in the report.

G.1.10.2. Differential Reports

Use this function to generate differential reports:

Postgres Pro Modules and Extensions Shipped as Individual Packages

```
get_diffreport([server name,] start1_id integer, end1_id integer, start2_id integer,
end2_id integer [, description text [, with_growth boolean [, db_exclude name[]]])
get_diffreport([server name,] time_range1 tstzrange, time_range2 tstzrange [, description
text [, with_growth boolean [, db_exclude name[]]])
get_diffreport([server name,] baseline1 varchar(25), baseline2 varchar(25) [, description
text [, with_growth boolean [, db_exclude name[]]])
get_diffreport([server name,] baseline1 varchar(25), time_range2 tstzrange [, description
text [, with_growth boolean[, db_exclude name[]]])
get_diffreport([server name,] time_range1 tstzrange, baseline2 varchar(25) [, description
text [, with_growth boolean[, db_exclude name[]]])
get_diffreport([server name,] start1_id integer, end1_id integer, baseline2 varchar(25)
[, description text [, with_growth boolean[, db_exclude name[]]])
get_diffreport([server name,] baseline1 varchar(25), start2_id integer, end2_id integer
[, description text [, with_growth boolean[, db_exclude name[]]])
```

Generates a differential report for two intervals. The combinations of arguments provide possible ways to specify the two intervals.

Arguments:

- *server* — server name. `local` sever is assumed if omitted.
- *start1_id*, *end1_id* — identifiers of the starting and ending samples for the first interval.
- *start2_id*, *end2_id* — identifiers of the starting and ending samples for the second interval.
- *baseline1* — baseline name for the first interval.
- *baseline2* — baseline name for the second interval.
- *time_range1* — time range for the first interval.
- *time_range2* — time range for the second interval.
- *description* — short text to be included in the report as its description.
- *with_growth* — flag requesting interval expansion to the nearest bounds with data on relation growth available. The default value is `false`.
- *db_exclude* — the database exclusion list. Lists databases to be excluded from all report tables having the `Database` column. Use to hide some databases in the report.

G.1.10.3. Report Generation Example

Generate a report for the `local` server and interval defined by samples:

```
psql -Aqtc "SELECT profile.get_report(480,482)" -o report_480_482.html
```

For any other server, provide its name:

```
psql -Aqtc "SELECT profile.get_report('omega',12,14)" -o report_omega_12_14.html
```

Generate a report using time ranges:

```
psql -Aqtc "SELECT profile.get_report(tstzrange('2020-05-13 11:51:35+03','2020-05-13
11:52:18+03'))" -o report_range.html
```

Generate a relative time-range report:

```
psql -Aqtc "SELECT profile.get_report(tstzrange(now() - interval '1 day',now()))" -o
report_last_day.html
```

G.1.11. pgpro_pwr Report Sections

Each `pgpro_pwr` report is divided into sections, described below. The number of top objects reported in each sorted report table is specified by the `pgpro_pwr.max` parameter.

Almost every item in the report can be accentuated by a single mouse click. The accentuated item will be instantly highlighted in all report sections, making it easy to find. The attributes identifying the item will appear in the bottom-right corner of the page. For example, if you click on a database name in the “Database statistics” report table, you can notice a small table with the database attributes in the bottom-right corner of the page.

When you scroll down the report, its table of contents will be available on the right side of the page. It can be hidden with a single mouse click on the “content” tag.

A substring-based filter is also available that helps limit the report contents to particular objects based on a substring. Specifically, substring-based filtering is applied to query texts.

G.1.11.1. Server statistics

Tables in this section of a `pgpro_pwr` report are described below.

The report table “Database statistics” provides per-database statistics for the report interval. The statistics are based on the `pg_stat_database` view. [Table G.1](#) lists columns of this report table.

Table G.1. Database statistics

Column	Description	Field/Calculation
Database	Database name	<code>datname</code>
Commits	Number of committed transactions	<code>xact_commit</code>
Rollbacks	Number of rolled back transactions	<code>xact_rollback</code>
Deadlocks	Number of deadlocks detected	<code>deadlocks</code>
Checksum Failures	Number of data page checksum failures detected in this database. This field is only shown if any checksum failures were detected in this database during the report interval.	<code>checksum_failures</code>
Checksums Last	Time at which the last data page checksum failure was detected in this database. This field is only shown if any checksum failures were detected in this database during the report interval.	<code>checksum_last_failure</code>
Hit%	Buffer cache hit ratio, i.e., percentage of pages fetched from buffers in all pages fetched	
Read	Number of disk blocks read in this database	<code>blks_read</code>
Hit	Number of times disk blocks were found already in the buffer cache	<code>blks_hit</code>
Ret	Number of returned tuples	<code>tup_returned</code>
Fet	Number of fetched tuples	<code>tup_fetched</code>
Ins	Number of inserted tuples	<code>tup_inserted</code>
Upd	Number of updated tuples	<code>tup_updated</code>
Del	Number of deleted tuples	<code>tup_deleted</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Temp Size	Total amount of data written to temporary files by queries in this database	temp_bytes
Temp Files	Number of temporary files created by queries in this database	temp_files
Size	Database size at the time of the last sample in the report interval	pg_database_size()
Growth	Database growth during the report interval	pg_database_size() increment between interval bounds

The report table “Cluster I/O statistics” provides I/O statistics by object types, backend types and contexts. This table is based on the `pg_stat_io` view of the [Cumulative Statistics System](#), available since Postgres Pro 16. [Table G.2](#) lists columns of this report table. Times are provided in seconds.

Table G.2. Cluster I/O statistics

Column	Description
Object	Target object of an I/O operation
Backend	Type of the backend that performed an I/O operation
Context	The context of an I/O operation
Reads Count	Number of read operations
Reads Bytes	Amount of data read
Reads Time	Time spent in read operations
Writes Count	Number of write operations
Writes Bytes	Amount of data written
Writes Time	Time spent in write operations
Writebacks Count	Number of blocks which the process requested the kernel write out to permanent storage
Writebacks Bytes	Amount of data requested for write out to permanent storage
Writebacks Time	Time spent in writeback operations, including the time spent queueing write-out requests and, potentially, the time spent to write out the dirty data
Extends Count	Number of relation extend operations
Extends Bytes	Amount of space used by extend operations
Extends Time	Time spent in extend operations
Hits	The number of times a desired block was found in a shared buffer
Evictions	Number of times a block has been written out from a shared or local buffer in order to make it available for another use
Reuses	The number of times an existing buffer in a size-limited ring buffer outside of shared buffers was reused as part of an I/O operation in the <code>bulkread</code> , <code>bulkwrite</code> , or <code>vacuum</code> contexts
Fsyncs Count	Number of fsync calls. These are only tracked in context <code>normal</code> .

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description
Fsyncs Time	Time spent in fsync operations

The report table “Cluster SLRU statistics” provides access statistics on SLRU (simple least-recently-used) caches. This table is based on the `pg_stat_slru` view of the [Cumulative Statistics System](#). [Table G.3](#) lists columns of this report table. Times are provided in seconds.

Table G.3. Cluster SLRU statistics

Column	Description	Field/Calculation
Name	Name of the SLRU	<code>name</code>
Zeroed	Number of blocks zeroed during initializations	<code>blks_zeroed</code>
Hits	Number of times disk blocks were found already in the SLRU, so that a read was not necessary (this only includes hits in the SLRU, not the operating system's file system cache)	<code>blks_hit</code>
Reads	Number of disk blocks read for this SLRU	<code>blks_read</code>
%Hit	Number of disk block hits for this SLRU as the percentage of $\text{Reads} + \text{Hits}$	$\text{blks_hit} * 100 / \text{blks_read} + \text{blks_hit}$
Writes	Number of disk blocks written for this SLRU	<code>blks_written</code>
Checked	Number of blocks checked for existence for this SLRU	<code>blks_exists</code>
Flushes	Number of flushes of dirty data for this SLRU	<code>flushes</code>
Truncates	Number of truncates for this SLRU	<code>truncates</code>

Table “Session statistics by database” is available in the report for Postgres Pro databases starting with version 14. This table is based on the `pg_stat_database` view of the [Statistics Collector](#). [Table G.4](#) lists columns of this report table. Times are provided in seconds.

Table G.4. Session statistics by database

Column	Description	Field/Calculation
Database	Database name	
Timing Total	Time spent by database sessions in this database during the report interval (note that statistics are only updated when the state of a session changes, so if sessions have been idle for a long time, this idle time won't be included)	<code>session_time</code>
Timing Active	Time spent executing SQL statements in this database during the report interval (this corresponds to the states <code>active</code> and <code>fastpath</code> function call in <code>pg_stat_activity</code>)	<code>active_time</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Timing Idle (T)	Time spent idling while in a transaction in this database during the report interval (this corresponds to the states <code>idle in transaction</code> and <code>idle in transaction (aborted)</code> in <code>pg_stat_activity</code>)	<code>idle_in_transaction_time</code>
Sessions Established	Total number of sessions established to this database during the report interval	<code>sessions</code>
Sessions Abandoned	Number of database sessions to this database that were terminated because connection to the client was lost during the report interval	<code>sessions_abandoned</code>
Sessions Fatal	Number of database sessions to this database that were terminated by fatal errors during the report interval	<code>sessions_fatal</code>
Sessions Killed	Number of database sessions to this database that were terminated by operator intervention during the report interval	<code>sessions_killed</code>

In Postgres Pro Shardman databases of versions that include [pgpro_stats](#) version starting with 1.4, workload statistics of vacuum processes are available. The “Database vacuum statistics” report table provides per-database aggregated total vacuum statistics based on the `pgpro_stats_vacuum_tables` view. [Table G.5](#) lists columns of this report table. Times are provided in seconds.

Table G.5. Database vacuum statistics

Column	Description	Field/Calculation
Database	Database name	
Blocks fetched	Total number of database blocks fetched by vacuum operations	<code>total_blks_read</code> + <code>total_blks_hit</code>
%Total	Total number of database blocks fetched (read+hit) by vacuum operations as the percentage of all blocks fetched in the cluster	<code>Blocks fetched * 100 / Cluster fetched</code>
Blocks read	Total number of database blocks read by vacuum operations	<code>total_blks_read</code>
%Total	Total number of database blocks read by vacuum operations as the percentage of all blocks read in the cluster	<code>Blocks read * 100 / Cluster read</code>
VM Frozen	Total number of blocks marked all-frozen in the visibility map	<code>pages_frozen</code>
VM Visible	Total number of blocks marked all-visible in the visibility map	<code>pages_all_visible</code>
Tuples deleted	Total number of dead tuples vacuum operations deleted from tables of this database	<code>tuples_deleted</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Tuples left	Total number of dead tuples vacuum operations left in tables of this database due to their visibility in transactions	dead_tuples
%Eff	Vacuum efficiency in terms of deleted tuples. This is the percentage of tuples deleted from tables of this database in all dead tuples to be deleted from tables of this database.	$\text{tuples_deleted} * 100 / (\text{tuples_deleted} + \text{dead_tuples})$
WAL size	Total amount of WAL bytes generated by vacuum operations performed on tables of this database	wal_bytes
Read I/O time	Time spent reading database blocks by vacuum operations performed on tables of this database	blk_read_time
Write I/O time	Time spent writing database blocks by vacuum operations performed on tables of this database	blk_write_time
%Total	Vacuum I/O time spent as the percentage of whole cluster I/O time	
Total vacuum time	Total time of vacuuming tables of this database	total_time
Delay vacuum time	Time spent sleeping in a vacuum delay point by vacuum operations performed on tables of this database	delay_time
User CPU time	User CPU time of vacuuming tables of this database	user_time
System CPU time	System CPU time of vacuuming tables of this database	system_time
Interrupts	Number of times vacuum operations performed on tables of this database were interrupted on any errors	interrupts

If the [pgpro_stats](#) extension supporting invalidation statistics was available during the report interval, the "Invalidation messages by database" report table provides per-database aggregated total invalidation message statistics. [Table G.6](#) lists columns of this report table.

Table G.6. Invalidation messages by database

Column	Description	Field/Calculation
Database	Database name	
Invalidation messages sent	Total number of invalidation messages sent by backends in this database. Statistics are provided for corresponding message	Fields of pgpro_stats_totals.inval_msgs

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	types of pgpro_stats_inval_msgs	
Cache resets	Total number of shared cache resets	pgpro_stats_totals .cache_resets

If the [pgpro_stats](#) extension was available during the report interval, the “Statement statistics by database” report table provides per-database aggregated total statistics for the `pgpro_stats_statements` view data. [Table G.7](#) lists columns of this report table. Times are provided in seconds.

Table G.7. Statement statistics by database

Column	Description	Field/Calculation
Database	Database name	
Calls	Number of times all statements in the database were executed	calls
Plan Time	Time spent planning all statements in the database	Sum of total_plan_time
Exec Time	Time spent executing all statements in the database	Sum of total_exec_time
Read Time	Time spent reading blocks by all statements in the database	Sum of blk_read_time
Write Time	Time spent writing blocks by all statements in the database	Sum of blk_write_time
Trg Time	Time spent executing trigger functions by all statements in the database	
Shared Fetched	Total number of shared blocks fetched by all statements in the database	Sum of (shared_blks_read + shared_blks_hit)
Local Fetched	Total number of local blocks fetched by all statements in the database	Sum of (local_blks_read + local_blks_hit)
Shared Dirtied	Total number of shared blocks dirtied by all statements in the database	Sum of shared_blks_dirtied
Local Dirtied	Total number of local blocks dirtied by all statements in the database	Sum of local_blks_dirtied
Read Temp	Total number of temp blocks read by all statements in the database	Sum of temp_blks_read
Write Temp	Total number of temp blocks written by all statements in the database	Sum of temp_blks_written
Read Local	Total number of local blocks read by all statements in the database	Sum of local_blks_read
Write Local	Total number of local blocks written by all statements in the database	Sum of local_blks_written

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Statements	Total number of captured statements	
WAL Size	Total amount of WAL generated by all statements in the database	Sum of wal_bytes

The “Statement average min/max timings” report table contains per-database aggregated min/max timing statistics from the one of [pgpro_stats](#) or [pg_stat_statements](#) extensions that was available during the report interval, with the precedence of [pgpro_stats](#). This report is sensitive to the fastest and slowest planning and execution of every statement in a cluster. That is, you can see execution and planning stability in your database. [Table G.8](#) lists columns of this report table.

Table G.8. Statement average min/max timings

Column	Description
Database	Database name
Min average planning time (ms)	The average value of min_plan_time for all statements and all samples included in the report
Max average planning time (ms)	The average value of max_plan_time for all statements and all samples included in the report
Delta% of average planning times	Difference between the mean max_plan_time and mean min_plan_time as the percentage of the mean min_plan_time. The less this difference, the more stable query planning in your database is.
Min average execution time (ms)	The average value of min_exec_time for all statements and all samples included in the report
Max average execution time (ms)	The average value of max_exec_time for all statements and all samples included in the report
Delta% of average execution times	Difference between the mean max_exec_time and mean min_exec_time as the percentage of the mean min_exec_time. The less this difference, the more stable query execution in your database is.
Statements	Total count of captured statements

If the JIT-related statistics was available in the statement statistics extension during the report interval, the “JIT statistics by database” report table provides per-database aggregated total statistics of JIT executions. [Table G.9](#) lists columns of this report table. Times are provided in seconds.

Table G.9. JIT statistics by database

Column	Description	Field/Calculation
Database	Database name	
Calls	Number of times all statements in the database were executed	calls
Plan Time	Time spent planning all statements in the database	Sum of total_plan_time
Exec Time	Time spent executing all statements in the database	Sum of total_exec_time
Generation count	Total number of functions JIT-compiled by the statements	Sum of jit_functions

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Generation time	Total time spent by the statements on generating JIT code	Sum of <code>jit_generation_time</code>
Inlining count	Number of times functions have been inlined	Sum of <code>jit_inlining_count</code>
Inlining time	Total time spent by statements on inlining functions	Sum of <code>jit_inlining_time</code>
Optimization count	Number of times statements have been optimized	Sum of <code>jit_optimization_count</code>
Optimization time	Total time spent by statements on optimizing	Sum of <code>jit_optimization_time</code>
Emission count	Number of times code has been emitted	Sum of <code>jit_emission_count</code>
Emission time	Total time spent by statements on emitting code	Sum of <code>jit_emission_time</code>
Deform count	Number of tuple deform functions JIT-compiled by the statement of the database	
Deform time	Total time spent by the statements of the database on JIT-compiling the tuple deform functions	

The report table “Cluster statistics” provides data from the `pg_stat_bgwriter` view. Table G.10 lists rows of this report table. Times are provided in seconds.

Table G.10. Cluster statistics

Row	Description	Field/Calculation
Scheduled checkpoints	Number of scheduled checkpoints that have been performed	<code>checkpoints_timed</code>
Requested checkpoints	Number of requested checkpoints that have been performed	<code>checkpoints_req</code>
Checkpoint write time (s)	Total amount of time that has been spent in the portion of checkpoint processing where files are written to disk	<code>checkpoint_write_time</code>
Checkpoint sync time (s)	Total amount of time that has been spent in the portion of checkpoint processing where files are synchronized to disk	<code>checkpoint_sync_time</code>
Checkpoint buffers written	Number of buffers written during checkpoints	<code>buffers_checkpoint</code>
Background buffers written	Number of buffers written by the background writer	<code>buffers_clean</code>
Backend buffers written	Number of buffers written directly by a backend. Will not be shown since Postgres Pro 17.	<code>buffers_backend</code>
Backend fsync count	Number of times a backend had to execute its own <code>fsync</code> call (normally the background writer handles those even when the	<code>buffers_backend_fsync</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Row	Description	Field/Calculation
	backend does its own write). Will not be shown since Postgres Pro 17.	
Bgwriter interrupts (too many buffers)	Number of times the background writer stopped a cleaning scan because it had written too many buffers	maxwritten_clean
Number of buffers allocated	Total number of buffers allocated	buffers_alloc
WAL generated	Total amount of WAL generated	pg_current_wal_lsn() value increment
Start LSN	Log sequence number at the start of a report interval	pg_current_wal_lsn() at the first sample of a report
End LSN	Log sequence number at the end of a report interval	pg_current_wal_lsn() at the last sample of a report
WAL generated by vacuum	Total amount of WAL generated by vacuum	Based on the wal_bytes field of the pgpro_stats_vacuum_databases view.
WAL segments archived	Total number of archived WAL segments	Based on pg_stat_archiver.archived_count
WAL segments archive failed	Total number of WAL segment archiver failures	Based on pg_stat_archiver.failed_count .
Archiver performance	Average archiver process performance per second	Based on the active_time field of the pgpro_stats_archiver view.
Archive command performance	Average archive_command performance per second	Based on the archive_command_time field of the pgpro_stats_archiver view.

Table “WAL statistics” is available in the report for Postgres Pro databases starting with version 14. This table is based on the [pg_stat_wal](#) view of the [Statistics Collector](#). [Table G.11](#) lists columns of this report table. Times are provided in seconds.

Table G.11. WAL statistics

Row	Description	Field/Calculation
WAL generated	Total amount of WAL generated during the report interval	wal_bytes
WAL per second	Average amount of WAL generated per second during the report interval	wal_bytes / report_duration
WAL records	Total number of WAL records generated during the report interval	wal_records
WAL FPI	Total number of WAL full page images generated during the report interval	wal_fpi
WAL buffers full	Number of times WAL data was written to disk because WAL buffers became full during the report interval	wal_buffers_full

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Row	Description	Field/Calculation
WAL writes	Number of times WAL buffers were written out to disk via <code>XLogWrite</code> request during the report interval	<code>wal_write</code>
WAL writes per second	Average number of times WAL buffers were written out to disk via <code>XLogWrite</code> request per second during the report interval	<code>wal_write / report_duration</code>
WAL sync	Number of times WAL files were synced to disk via <code>issue_xlog_fsync</code> request during the report interval (if <code>fsync</code> is on and <code>wal_sync_method</code> is either <code>fdatasync</code> , <code>fsync</code> or <code>fsync_writethrough</code> , otherwise zero). See Section 28.5 for more information about the internal WAL function <code>issue_xlog_fsync</code> .	<code>wal_sync</code>
WAL syncs per second	Average number of times WAL files were synced to disk via <code>issue_xlog_fsync</code> request per second during the report interval	<code>wal_sync / report_duration</code>
WAL write time	Total amount of time spent writing WAL buffers to disk via <code>XLogWrite</code> request during the report interval (if <code>track_wal_io_timing</code> is enabled, otherwise zero; for more details, see Section 19.9). This includes the sync time when <code>wal_sync_method</code> is either <code>open_datasync</code> or <code>open_sync</code> .	<code>wal_write_time</code>
WAL write duty	WAL write time as the percentage of the report interval duration	<code>wal_write_time * 100 / report_duration</code>
WAL sync time	Total amount of time spent syncing WAL files to disk via <code>issue_xlog_fsync</code> request during the report interval (if <code>track_wal_io_timing</code> is enabled, <code>fsync</code> is on, and <code>wal_sync_method</code> is either <code>fdatasync</code> , <code>fsync</code> or <code>fsync_writethrough</code> , otherwise zero).	<code>wal_sync_time</code>
WAL sync duty	WAL sync time as the percentage of the report interval duration	<code>wal_sync_time * 100 / report_duration</code>

The report table “Tablespace statistics” provides information on the sizes and growth of tablespaces. [Table G.12](#) lists columns of this report table.

Table G.12. Tablespace statistics

Column	Description	Field/Calculation
Tablespace	Tablespace name	<code>pg_tablespace .spcname</code>
Path	Tablespace path	<code>pg_tablespace_location()</code>
Size	Tablespace size at the time of the last sample in the report interval	<code>pg_tablespace_size()</code>
Growth	Tablespace growth during the report interval	<code>pg_tablespace_size()</code> increment between interval bounds

If the [pgpro_stats](#) extension was available during the report interval, the report table “Wait statistics by database” shows the total wait time by wait event type and database. [Table G.13](#) lists columns of this report table.

Table G.13. Wait statistics by database

Column	Description
Database	Database name
Wait event type	Type of event for which the backends were waiting. Asterisk means aggregation of all wait event types in the database.
Waited (s)	Time spent waiting in events of <code>Wait event type</code> , in seconds
%Total	Percentage of wait time spent in the database events of <code>Wait event type</code> in all wait time for the cluster

If the [pgpro_stats](#) extension was available during the report interval, the report table “Top wait events” shows top wait events in the cluster by wait time. [Table G.14](#) lists columns of this report table.

Table G.14. Top wait events

Column	Description
Database	Database name
Wait event type	The type of event for which the backends were waiting
Wait event	Wait event name for which the backends were waiting
Waited	Total wait time spent in <code>Wait event</code> of the database, in seconds
%Total	Percentage of wait time spent in <code>Wait event</code> of the database in all wait time in the cluster

G.1.11.2. Load distribution

This section of a `pgpro_pwr` report is based on the `pgpro_stats_totals` view of the [pgpro_stats](#) extension if it was available during the report interval. Each table in this section provides data for the report interval on load distribution for a certain kind of objects for which aggregated statistics are collected, such as databases, applications, hosts, or users. Each table contains one row for each resource (for example, total time or shared blocks written), where load distribution is shown in graphics, as a stacked bar chart for top objects by load of this resource. If the bar chart area that corresponds to an object is too narrow to include captions, point that area to get a hint with the caption, value and percentage.

The report tables “Load distribution among heavily loaded databases”, “Load distribution among heavily loaded applications”, “Load distribution among heavily loaded hosts” and “Load distribution among heavily loaded users” show load distribution for respective objects. [Table G.15](#) lists rows of these report tables.

Table G.15. Load distribution

Row	Description	Calculation
Total time (sec.)	Total time spent in the planning and execution of statements	<code>total_plan_time + total_exec_time</code>
Executed count	Number of queries executed	<code>queries_executed</code>
I/O time (sec.)	Total time the statements spent reading or writing blocks (if track_io_timing is enabled, otherwise zero)	<code>blk_read_time + blk_write_time</code>
Blocks fetched	Total number of shared block cache hits and shared blocks read by the statements	<code>shared_blks_hit + shared_blks_read</code>
Shared blocks read	Total number of shared blocks read by the statements	<code>shared_blks_read</code>
Shared blocks dirtied	Total number of shared blocks dirtied by the statements	<code>shared_blks_dirtied</code>
Shared blocks written	Total number of shared blocks written by the statements	<code>shared_blks_written</code>
WAL generated	Total amount of WAL generated by the statements	<code>wal_bytes</code>
Temp and Local blocks written	Total number of temporary and local blocks written by the statements	<code>temp_blks_written + local_blks_written</code>
Temp and Local blocks read	Total number of temp and local blocks read by the statements	<code>temp_blks_read + local_blks_read</code>
Invalidation messages sent	Total number of all invalidation messages sent by backends in this database	<code>(pgpro_stats_totals.inval_msgs).all</code>
Cache resets	Total number of shared cache resets	<code>pgpro_stats_totals.cache_resets</code>

G.1.11.3. Session states observed by subsamples

This section of a `pgpro_pwr` report provides information about session states captured by subsamples during the report interval.

Tables of this report section are described below.

The report subsection “Chart with session state” visualizes session states captured by subsamples. It is the timeline chart showing captured session states in backends and transactions. Every state contains a popup with session state attributes. Click on a state to see this state in the table of session states.

The report table “Session state statistics by database” shows the aggregated data on session states. Only session states captured in subsamples are counted. [Table G.16](#) lists columns of this report table.

Table G.16. Session state statistics by database

Column	Description
Database	Database name

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description
Summary Active	Overall time of active states captured in subsamples
Summary Idle in xact	Overall time of idle in transaction states captured in subsamples
Summary Idle in xact (A)	Overall time of idle in transaction (aborted) states captured in subsamples
Maximal Active	Time of the longest active state captured in subsamples
Maximal Idle in xact	Time of the longest idle in transaction state captured in subsamples
Maximal Idle in xact (A)	Time of the longest idle in transaction (aborted) state captured in subsamples
Maximal xact age	Maximal transaction age detected in subsamples

The report table “Top 'idle in transaction' session states by duration” shows top `pgpro_pwr.topn` longest idle in transaction states that were last observed in the [pg_stat_activity](#) view for each session. [Table G.17](#) lists columns of this report table.

Table G.17. Top 'idle in transaction' session states by duration

Column	Description	Field/Calculation
Database	Database name	<code>datname</code>
User	User name	<code>username</code>
App	Application name	<code>application_name</code>
Pid	Process ID	<code>pid</code>
Xact start	Transaction start timestamp	<code>xact_start</code>
State change	State change timestamp	<code>state_change</code>
State duration	State duration	<code>clock_timestamp() - state_change</code>

The report table “Top 'active' session states by duration” shows top `pgpro_pwr.topn` longest active states that were last observed in the [pg_stat_activity](#) view for each session. [Table G.18](#) lists columns of this report table.

Table G.18. Top 'active' session states by duration

Column	Description	Field/Calculation
Database	Database name	<code>datname</code>
User	User name	<code>username</code>
App	Application name	<code>application_name</code>
Pid	Process ID	<code>pid</code>
Xact start	Transaction start timestamp	<code>xact_start</code>
State change	State change timestamp	<code>state_change</code>
State duration	State duration	<code>clock_timestamp() - state_change</code>

The report table “Top states by transaction age” shows top session states by transaction age that were last observed in the [pg_stat_activity](#) view for each session. [Table G.19](#) lists columns of this report table.

Table G.19. Top states by transaction age

Column	Description	Field/Calculation
Database	Database name	datname
User	User name	username
App	Application name	application_name
Pid	Process ID	pid
Xact start	Transaction start timestamp	xact_start
Xact duration	Transaction duration	clock_timestamp() - xact_start
Age	Transaction age	age(backend_xmin)
State	Session state at the maximum age detected	
State change	State change timestamp	state_change
State duration	State duration	clock_timestamp() - state_change

The report table “Top states by transaction duration” shows top longest session states that were last observed in the [pg_stat_activity](#) view for each session. [Table G.20](#) lists columns of this report table.

Table G.20. Top states by transaction duration

Column	Description	Field/Calculation
Database	Database name	datname
User	User name	username
App	Application name	application_name
Pid	Process ID	pid
Xact start	Transaction start timestamp	xact_start
Xact duration	Transaction duration	clock_timestamp() - xact_start
Age	Transaction age	age(backend_xmin)
State	Session state at the maximum age detected	
State change	State change timestamp	state_change
State duration	State duration	clock_timestamp() - state_change

G.1.11.4. SQL query statistics

This section of a pgpro_pwr report provides data for the report interval on top statements by several important statistics. The data is mainly captured from views of the one of [pgpro_stats](#) and [pg_stat_statements](#) extensions that was available during the report interval, with the precedence of pgpro_stats. Each statement can be highlighted in all SQL-related sections with a single mouse click on it. This click will also show a query text preview just under the query statistics row. The query text preview can be hidden with a second click on a query.

Tables of this report section are described below.

The report table “Top SQL by elapsed time” shows top statements by the sum of `total_plan_time` and `total_exec_time` fields of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.21](#) lists columns of this report table. Times are provided in seconds.

Table G.21. Top SQL by elapsed time

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
%Total	Percentage of elapsed time of this statement plan in the total elapsed time of all statements in the cluster	
Elapsed Time (s)	Total time spent in planning and execution of the statement plan	<code>total_plan_time</code> + <code>total_exec_time</code>
Plan Time (s)	Total time spent in planning of the statement	<code>total_plan_time</code>
Exec Time (s)	Total time spent in execution of the statement plan	<code>total_exec_time</code>
JIT Time (s)	Total time spent by JIT executing this statement plan, in seconds	<code>jit_generation_time</code> + <code>jit_inlining_time</code> + <code>jit_optimization_time</code> + <code>jit_emission_time</code>
Read I/O time (s)	Total time the statement spent reading blocks	<code>blk_read_time</code>
Write I/O time (s)	Total time the statement spent writing blocks	<code>blk_write_time</code>
Usr CPU time (s)	Time spent on CPU in the user space, in seconds	<code>rusage.user_time</code>
Sys CPU time (s)	Time spent on CPU in the system space, in seconds	<code>rusage.system_time</code>
Plans	Number of times the statement was planned	<code>plans</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by planning time” shows top statements by the value of the `total_plan_time` field of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.22](#) lists columns of this report table.

Table G.22. Top SQL by planning time

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Plan elapsed(s)	Total time spent in planning of the statement, in seconds	<code>total_plan_time</code>
%Elapsed	Percentage of <code>total_plan_time</code> in the sum of <code>total_plan_time</code> and <code>total_exec_time</code> of this statement plan	
Mean plan time	Mean time spent planning the statement, in milliseconds	<code>mean_plan_time</code>
Min plan time	Minimum time spent planning the statement, in milliseconds	<code>min_plan_time</code>
Max plan time	Maximum time spent planning the statement, in milliseconds	<code>max_plan_time</code>
StdErr plan time	Population standard deviation of time spent planning the statement, in milliseconds	<code>stddev_plan_time</code>
Plans	Number of times the statement was planned	<code>plans</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by execution time” shows top statements by the value of the `total_time` field of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.23](#) lists columns of this report table.

Table G.23. Top SQL by execution time

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Exec (s)	Total time spent executing the statement plan, in seconds	total_exec_time
%Elapsed	Percentage of total_exec_time of this statement plan in this statement elapsed time	
%Total	Percentage of total_exec_time of this statement plan in the total elapsed time of all statements in the cluster	
JIT Time (s)	Total time spent by JIT executing this statement plan, in seconds	jit_generation_time + jit_inlining_time + jit_optimization_time + jit_emission_time
Read I/O time (s)	Total time spent in reading pages while executing the statement plan, in seconds	blk_read_time
Write I/O time (s)	Total time spent in writing pages while executing the statement plan, in seconds	blk_write_time
Usr CPU time (s)	Time spent on CPU in the user space, in seconds	rusage.user_time
Sys CPU time (s)	Time spent on CPU in the system space, in seconds	rusage.system_time
Rows	Number of rows retrieved or affected by execution of the statement plan	rows
Mean execution time	Mean time spent executing the statement plan, in milliseconds	mean_exec_time
Min execution time	Minimum time spent executing the statement plan, in milliseconds	min_exec_time
Max execution time	Maximum time spent executing the statement plan, in milliseconds	max_exec_time
StdErr execution time	Population standard deviation of time spent executing the statement plan, in milliseconds	stddev_exec_time
Executions	Number of executions of this statement plan	calls
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

The report table “Top SQL by mean execution time” shows top `pgpro_pwr.topn` statements by the value of the `mean_time` or `mean_exec_time` field of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.24](#) lists columns of this report table.

Table G.24. Top SQL by mean execution time

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Mean execution time (ms)	Mean time spent executing the statement, in milliseconds	<code>mean_exec_time</code>
Min execution time (ms)	Minimum time spent executing the statement, in milliseconds	<code>min_exec_time</code>
Max execution time (ms)	Maximum time spent executing the statement, in milliseconds	<code>max_exec_time</code>
StdErr execution time	Population standard deviation of time spent executing the statement, in milliseconds	<code>stddev_exec_time (ms)</code>
Exec (s)	Time spent executing this statement, in seconds	<code>total_exec_time</code>
%Elapsed	Execution time of this statement as the percentage of the statement elapsed time	
%Total	Execution time of this statement as the percentage of the total elapsed time of all statements in the cluster	
JIT time (s)	Total time spent by JIT executing this statement, in seconds	<code>jit_generation_time + jit_inlining_time + jit_optimization_time + jit_emission_time</code>
Read I/O time (s)	Time spent reading blocks, in seconds	<code>blk_read_time</code>
Write I/O time (s)	Time spent writing blocks, in seconds	<code>blk_write_time</code>
Rows	Number of rows retrieved or affected by the statement	<code>rows</code>
Executions	Number of times this statement was executed	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by executions” shows top statements by the value of the `calls` field of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.25](#) lists columns of this report table.

Table G.25. Top SQL by executions

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Total	Percentage of <code>calls</code> of this statement plan in the total <code>calls</code> of all statements in the cluster	
Rows	Number of rows retrieved or affected by execution of the statement plan	<code>rows</code>
Mean (ms)	Mean time spent executing the statement plan, in milliseconds	<code>mean_exec_time</code>
Min (ms)	Minimum time spent executing the statement plan, in milliseconds	<code>min_exec_time</code>
Max (ms)	Maximum time spent executing the statement plan, in milliseconds	<code>max_exec_time</code>
StdErr (ms)	Population standard deviation of time spent executing the statement plan, in milliseconds	<code>stddev_time</code>
Elapsed(s)	Total time spent executing the statement plan, in seconds	<code>total_exec_time</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by I/O wait time” shows top statements by read and write time, i.e., sum of values of `blk_read_time` and `blk_write_time` fields of the `pgpro_stats_statements` or `pg_stat_statements` view. [Table G.26](#) lists columns of this report table. Times are provided in seconds.

Table G.26. Top SQL by I/O wait time

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	(such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
IO(s)	Total time spent in reading and writing while executing this statement plan, i.e., I/O time	blk_read_time + blk_write_time
R(s)	Total time spent in reading while executing this statement plan	blk_read_time
W(s)	Total time spent in writing while executing this statement plan	blk_write_time
%Total	Percentage of I/O time of this statement plan in the total I/O time of all statements in the cluster	
Shr Reads	Total number of shared blocks read while executing the statement plan	shared_blks_read
Loc Reads	Total number of local blocks read while executing the statement plan	local_blks_read
Tmp Reads	Total number of temp blocks read while executing the statement plan	temp_blks_read
Shr Writes	Total number of shared blocks written while executing the statement plan	shared_blks_written
Loc Writes	Total number of local blocks written while executing the statement plan	local_blks_written
Tmp Writes	Total number of temp blocks written while executing the statement plan	temp_blks_written
Elapsed(s)	Total time spent in execution of the statement plan	total_plan_time + total_exec_time
Executions	Number of executions of the statement plan	calls
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by shared blocks fetched” shows top statements by the number of read and hit blocks, which helps to detect the most data-intensive statements. [Table G.27](#) lists columns of this report table.

Table G.27. Top SQL by shared blocks fetched

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Blks fetched	Number of blocks retrieved while executing the statement plan	<code>shared_blks_hit</code> + <code>shared_blks_read</code>
%Total	Percentage of blocks fetched while executing the statement plan in all blocks fetched for all statements in the cluster	
Hits(%)	Percentage of blocks got from buffers in all blocks got	
Elapsed(s)	Total time spent in execution of the statement plan, in seconds	<code>total_plan_time</code> + <code>total_exec_time</code>
Rows	Number of rows retrieved or affected by execution of the statement plan	<code>rows</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by shared blocks read” shows top statements by the number of shared reads, which helps to detect the most read-intensive statements. [Table G.28](#) lists columns of this report table.

Table G.28. Top SQL by shared blocks read

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
User	Name of the user executing the statement	Derived from <code>userid</code>
Reads	Number of shared blocks read while executing this statement plan	<code>shared_blks_read</code>
%Total	Percentage of shared reads for this statement plan in all shared reads of all statements in the cluster	
Hits (%)	Percentage of blocks got from buffers in all blocks got while executing this statement plan	
Elapsed (s)	Total time spent in execution of the statement plan, in seconds	<code>total_plan_time</code> + <code>total_exec_time</code>
Rows	Number of rows retrieved or affected by execution of the statement plan	<code>rows</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by shared blocks dirtied” shows top statements by the number of shared dirtied buffers, which helps to detect statements that do most data changes in the cluster. [Table G.29](#) lists columns of this report table.

Table G.29. Top SQL by shared blocks dirtied

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Dirtied	Number of shared buffers dirtied while executing this statement plan	<code>shared_blks_dirtied</code>
%Total	Percentage of dirtied shared buffers for this statement plan in all dirtied shared buffers of all statements in the cluster	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Hits (%)	Percentage of blocks got from buffers in all blocks got while executing this statement plan	
WAL	Total amount of WAL bytes generated by the statement plan	wal_bytes
%Total	Percentage of WAL bytes generated by the statement plan in total WAL generated in the cluster	
Elapsed(s)	Total time spent in execution of the statement plan, in seconds	total_plan_time + total_exec_time
Rows	Number of rows retrieved or affected by execution of the statement plan	rows
Executions	Number of executions of the statement plan	calls
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by shared blocks written” shows top statements by the number of blocks written. [Table G.30](#) lists columns of this report table.

Table G.30. Top SQL by shared blocks written

Column	Description	Field/Calculation
Query ID	Hex representation of queryid. The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Written	Number of blocks written while executing this statement plan	shared_blks_written
%Total	Percentage of blocks written by this statement plan in all written blocks in the cluster	Percentage of shared_blks_written in (pg_stat_bgwriter.buffer_checkpoint + pg_stat_bgwriter.buffer_clean + pg_stat_bgwriter.buffer_backend)
%BackendW	Percentage of blocks written by this statement plan in all blocks in the cluster written by backends	Percentage of shared_blks_written in pg_stat_bgwriter.buffer_backend

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Hits (%)	Percentage of blocks got from buffers in all blocks got while executing this statement plan	
Elapsed(s)	Total time spent in execution of the statement plan, in seconds	<code>total_plan_time</code> + <code>total_exec_time</code>
Rows	Number of rows retrieved or affected by execution of the statement plan	<code>rows</code>
Executions	Number of executions of the statement plan	<code>calls</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by WAL size” shows top statements by the amount of WAL generated. [Table G.31](#) lists columns of this report table.

Table G.31. Top SQL by WAL size

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
WAL	Total amount of WAL bytes generated by the statement plan	<code>wal_bytes</code>
%Total	Percentage of WAL bytes generated by the statement plan in total WAL generated in the cluster	
Dirtied	Number of shared buffers dirtied while executing this statement plan	<code>shared_blks_dirtied</code>
WAL FPI	Total number of WAL full page images generated by the statement plan	<code>wal_fpi</code>
WAL records	Total number of WAL records generated by the statement plan	<code>wal_records</code>
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

The report table “Top SQL by temp usage” shows top statements by temporary I/O, which is calculated as the sum of `temp_blks_read`, `temp_blks_written`, `local_blks_read` and `local_blks_written` fields. Table G.32 lists columns of this report table.

Table G.32. Top SQL by temp usage

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	<code>planid</code>
Database	Database name for the statement	Derived from <code>dbid</code>
User	Name of the user executing the statement	Derived from <code>userid</code>
Local fetched	Number of local blocks retrieved	<code>local_blks_hit</code> + <code>local_blks_read</code>
Hits(%)	Percentage of local blocks got from buffers in all local blocks got	
Write Local (blk)	Number of blocks written by this statement plan that are used in temporary tables	<code>local_blks_written</code>
Write Local %Total	Percentage of <code>local_blks_written</code> of this statement plan in the total of <code>local_blks_written</code> for all statements in the cluster	
Read Local (blk)	Number of blocks read by this statement plan that are used in temporary tables	<code>local_blks_read</code>
Read Local %Total	Percentage of <code>local_blks_read</code> of this statement plan in the total of <code>local_blks_read</code> for all statements in the cluster	
Write Temp (blk)	Number of temp blocks written by this statement plan	<code>temp_blks_written</code>
Write Temp %Total	Percentage of <code>temp_blks_written</code> of this statement plan in the total of <code>temp_blks_written</code> for all statements in the cluster	
Read Temp (blk)	Number of temp blocks read by this statement plan	<code>temp_blks_read</code>
Read Temp %Total	Percentage of <code>temp_blks_read</code> of this statement plan in the to-	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	total of temp_blks_read for all statements in the cluster	
Elapsed(s)	Total time spent in execution of the statement plan, in seconds	total_plan_time + total_exec_time
Rows	Number of rows retrieved or affected by execution of the statement plan	rows
Executions	Number of executions of the statement plan	calls
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

The report table “Top SQL by invalidation messages sent” shows top statements by the number of invalidation messages sent. [Table G.33](#) lists columns of this report table.

Table G.33. Top SQL by invalidation messages sent

Column	Description	Field/Calculation
Query ID	Hex representation of queryid. The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	queryid
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Invalidation messages sent	Total number of invalidation messages sent by backends executing this statement. Statistics are provided for corresponding message types of pgpro_stats_inval_msgs	fields of pgpro_stats_statements.inval_msgs

G.1.11.4.1. rusage statistics

This section is included in the report only if the [pgpro_stats](#) or pg_stat_kcache extension was available during the report interval.

The report table “Top SQL by system and user time” shows top statements by the sum of user_time and system_time fields of pg_stat_kcache or of the pgpro_stats_totals view. [Table G.34](#) lists columns of this report table.

Table G.34. Top SQL by system and user time

Column	Description	Field/Calculation
Query ID	Hex representation of queryid. The hash of the query ID, database ID and user ID is in square	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Plan User (s)	User CPU time elapsed during planning, in seconds	plan_user_time
Exec User (s)	User CPU time elapsed during execution, in seconds	exec_user_time
User %Total	Percentage of <code>plan_user_time + exec_user_time</code> in the total user CPU time for all statements	
Plan System (s)	System CPU time elapsed during planning, in seconds	plan_system_time
Exec System (s)	System CPU time elapsed during execution, in seconds	exec_system_time
System %Total	Percentage of <code>plan_system_time + exec_system_time</code> in the total system CPU time for all statements	

The report table “Top SQL by reads/writes done by filesystem layer” shows top statements by the sum of reads and writes fields of `pg_stat_kcache`. [Table G.35](#) lists columns of this report table.

Table G.35. Top SQL by reads/writes done by filesystem layer

Column	Description	Field/Calculation
Query ID	Hex representation of <code>queryid</code> . The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Plan Read Bytes	Bytes read during planning	plan_reads
Exec Read Bytes	Bytes read during execution	exec_reads
Read Bytes %Total	Percentage of <code>plan_reads + exec_reads</code> in the total number of bytes read by the filesystem layer for all statements	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Plan Writes	Bytes written during planning	plan_writes
Exec Writes	Bytes written during execution	exec_writes
Write %Total	Percentage of plan_writes + exec_writes in the total number of bytes written by the filesystem layer for all statements	

G.1.11.5. SQL query wait statistics

If the [pgpro_stats](#) extension was available during the report interval, this section of the report will contain a table that is split into sections, each showing top statements by overall wait time or by wait time for a certain [wait event type](#). Table sections related to specific wait events follow in the descending order of the total wait time in wait events of this type. [Table G.36](#) lists columns of this report table. Times are provided in seconds.

Table G.36. SQL query wait statistics

Column	Description	Field/Calculation
Query ID	Hex representation of queryid. The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	
Plan ID	Hash of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
Waited	Total wait time for all wait events of this statement plan	
%Total	Percentage of the total wait time of this statement plan in all the wait time in the cluster	
Details	Waits of this statement plan by wait types	

If the JIT-related statistics was available in the statement statistics extension during the report interval, the report table “Top SQL by JIT elapsed time” shows top statements by the sum of jit_*_time fields of the pgpro_stats_statements or pg_stat_statements view. [Table G.37](#) lists columns of this report table. Times are provided in seconds.

Table G.37. Top SQL by JIT elapsed time

Column	Description	Field/Calculation
Query ID	Hex representation of queryid. The hash of the query ID, database ID and user ID is in square brackets. The (N) mark will appear here for nested statements (such as statements invoked within top-level statements).	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Plan ID	Internal hash code, computed from the tree of the statement plan	planid
Database	Database name for the statement	Derived from dbid
User	Name of the user executing the statement	Derived from userid
JIT Time (s)	Total time spent by JIT executing this statement plan	jit_generation_time + jit_inlining_time + jit_optimization_time + jit_emission_time
Generation count	Total number of functions JIT-compiled by this statement	Sum of jit_functions
Generation time	Total time spent by this statement on generating JIT code	Sum of jit_generation_time
Inlining count	Number of times functions have been inlined	Sum of jit_inlining_count
Inlining time	Total time spent by this statement on inlining functions	Sum of jit_inlining_time
Optimization count	Number of times this statement has been optimized	Sum of jit_optimization_count
Optimization time	Total time spent by this statement on optimizing	Sum of jit_optimization_time
Emission count	Number of times code has been emitted	Sum of jit_emission_count
Emission time	Total time spent by this statement on emitting code	Sum of jit_emission_time
Deform count	Number of tuple deform functions JIT-compiled by the statement	
Deform time	Total time spent by the statement on JIT-compiling the tuple deform	
Plan Time (s)	Total time spent in planning of the statement	total_plan_time
Exec Time (s)	Total time spent in execution of the statement plan	total_exec_time
Read I/O time (s)	Total time the statement spent reading blocks	blk_read_time
Write I/O time (s)	Total time the statement spent writing blocks	blk_write_time
%Cvr	Coverage: duration of statement statistics collection as the percentage of the report duration	

G.1.11.6. Complete list of SQL texts

The “Complete list of SQL texts” section of the report contains a table that provides query and plan texts for all statements mentioned in the report. Use an appropriate Query ID/Plan ID link in any statistics table to see the corresponding query/plan text. [Table G.38](#) lists columns of this report table.

Table G.38. Complete list of SQL texts

Column	Description
ID	Hex representation of the query or plan identifier
Query/Plan Text	Text of the query or statement plan

G.1.11.7. Schema object statistics

Tables in this section of the report show top database objects by statistics from the Postgres Pro's [Statistics Collector](#) views. Report tables that contain data for tables and indexes provide a preview of storage parameters. You can click on a row to see storage parameters of the object right under the row.

The report table “Top tables by estimated sequentially scanned volume” shows top tables by estimated volume read by sequential scans. This can help you find database tables that possibly lack some index. When there are no relation sizes collected with `pg_relation_size()`, relation-size estimates are based on the `pg_class.relpages` field. Since such values are less accurate, they are shown in square brackets. The data is based on the `pg_stat_all_tables` view. [Table G.39](#) lists columns of this report table.

Table G.39. Top tables by estimated sequentially scanned volume

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
~SeqBytes	Estimated volume read by sequential scans	Sum of (<code>pg_relation_size()</code> * <code>seq_scan</code>)
SeqScan	Number of sequential scans performed on the table	<code>seq_scan</code>
IxScan	Number of index scans initiated on the table	<code>idx_scan</code>
IxFet	Number of live rows fetched by index scans	<code>idx_tup_fetch</code>
Ins	Number of rows inserted	<code>n_tup_ins</code>
Upd	Number of rows updated	<code>n_tup_upd</code>
Del	Number of rows deleted	<code>n_tup_del</code>
Upd (HOT)	Number of rows HOT updated	<code>n_tup_hot_upd</code>

In the report table “Top tables by blocks fetched”, blocks fetched include blocks being processed from disk (read) and from shared buffers (hit). This report table shows top database tables by the sum of blocks fetched for the table's heap, indexes, TOAST table (if any) and TOAST table index (if any). This can help you focus on tables with excessive processing of blocks. The data is based on the `pg_statio_all_tables` view. [Table G.40](#) lists columns of this report table.

Table G.40. Top tables by blocks fetched

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Table	Table name	
Heap Blks	Number of blocks fetched for the table's heap	heap_blks_read + heap_blks_hit
Heap Blks %Total	Percentage of blocks fetched for the table's heap in all blocks fetched in the cluster	
Ix Blks	Number of blocks fetched for table's indexes	idx_blks_read + idx_blks_hit
Ix Blks %Total	Percentage of blocks fetched for table's indexes in all blocks fetched in the cluster	
TOAST Blks	Number of blocks fetched for the table's TOAST table	toast_blks_read + toast_blks_hit
TOAST Blks %Total	Percentage of blocks fetched for the table's TOAST table in all blocks fetched in the cluster	
TOAST-Ix Blks	Number of blocks fetched for the table's TOAST index	tidx_blks_read + tidx_blks_hit
TOAST-Ix Blks %Total	Percentage of blocks fetched for the table's TOAST index in all blocks fetched in the cluster	

The report table “Top tables by blocks read” shows top database tables by the number of blocks read for the table's heap, indexes, TOAST table (if any) and TOAST table index (if any). This can help you focus on tables with excessive block readings. The data is based on the [pg_statio_all_tables](#) view. [Table G.41](#) lists columns of this report table.

Table G.41. Top tables by blocks read

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Heap Blks	Number of blocks read for the table's heap	heap_blks_read
Heap Blks %Total	Percentage of blocks read from the table's heap in all blocks read in the cluster	
Ix Blks	Number of blocks read from table's indexes	idx_blks_read
Ix Blks %Total	Percentage of blocks read from table's indexes in all blocks read in the cluster	
TOAST Blks	Number of blocks read from the table's TOAST table	toast_blks_read
TOAST Blks %Total	Percentage of blocks read from the table's TOAST table in all blocks read in the cluster	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
TOAST-Ix Blks	Number of blocks read from the table's TOAST index	tidx_blks_read
TOAST-Ix Blks %Total	Percentage of blocks read from the table's TOAST index in all blocks read in the cluster	
Hit (%)	Percentage of table, index, TOAST and TOAST index blocks got from buffers for this table in all blocks got for this table from either file system or buffers	

The report table “Top DML tables” shows top tables by the number of DML-affected rows, i.e., by the sum of `n_tup_ins`, `n_tup_upd` and `n_tup_del` (including TOAST tables). The data is based on the `pg_stat_all_tables` view. Table G.42 lists columns of this report table.

Table G.42. Top DML tables

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Ins	Number of rows inserted	<code>n_tup_ins</code>
Upd	Number of rows updated, including HOT	<code>n_tup_upd</code>
Del	Number of rows deleted	<code>n_tup_del</code>
Upd (HOT)	Number of rows HOT updated	<code>n_tup_hot_upd</code>
SeqScan	Number of sequential scans performed on the table	<code>seq_scan</code>
SeqFet	Number of live rows fetched by sequential scans	<code>seq_tup_read</code>
IxScan	Number of index scans initiated on this table	<code>idx_scan</code>
IxFet	Number of live rows fetched by index scans	<code>idx_tup_fetch</code>

The report table “Top tables by updated/deleted tuples” shows top tables by tuples modified by UPDATE/DELETE operations, i.e., by the sum of `n_tup_upd` and `n_tup_del` (including TOAST tables). The data is based on the `pg_stat_all_tables` view. Table G.43 lists columns of this report table.

Table G.43. Top tables by updated/deleted tuples

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Upd	Number of rows updated, including HOT	<code>n_tup_upd</code>

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Upd (HOT)	Number of rows HOT updated	n_tup_hot_upd
Del	Number of rows deleted	n_tup_del
Vacuum	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	vacuum_count
AutoVacuum	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count
Analyze	Number of times this table was manually analyzed	analyze_count
AutoAnalyze	Number of times this table was analyzed by the autovacuum daemon	autoanalyze_count

The report table “Top tables by removed all-visible marks” shows top tables by the number of times that the all-visible mark was removed from the visibility map by any backend. This report section is only shown when corresponding statistics are available. [Table G.44](#) lists columns of this report table.

Table G.44. Top tables by removed all-visible marks

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
All-Visible marks cleared	Number of times that the all-visible mark was removed from the relation visibility map	rev_all_visible_pages
All-Visible marks set	Number of times that the all-visible mark was set in the relation visibility map	pages_all_visible
All-Visible marks %Set	Percentage of the number of times that the all-visible mark was set in the number of times that it was set or removed	$\text{pages_all_visible} * 100\% / (\text{rev_all_visible_pages} + \text{pages_all_visible})$
Vacuum	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	vacuum_count
AutoVacuum	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count

The report table “Top tables by removed all-frozen marks” shows top tables by the number of times that the all-frozen mark was removed from the visibility map by any backend. This report section is only shown when corresponding statistics are available. [Table G.45](#) lists columns of this report table.

Table G.45. Top tables by removed all-frozen marks

Column	Description	Field/Calculation
DB	Database name for the table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
All-Frozen marks cleared	Number of times that the all-frozen mark was removed from the relation visibility map	rev_all_frozen_pages
All-Frozen marks set	Number of times that the all-frozen mark was set in the relation visibility map	pages_frozen
All-Frozen marks %Set	Percentage of the number of times that the all-frozen mark was set in the number of times that it was set or removed	$\text{pages_frozen} * 100\% / (\text{rev_all_frozen_pages} + \text{pages_frozen})$
Vacuum	Number of times this table has been manually vacuumed (not counting VACUUM FULL)	vacuum_count
AutoVacuum	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count

The report table “Top tables by new-page updated tuples” shows top tables by the number of rows updated where the successor version goes onto a *new* heap page, leaving behind an original version with a `t_ctid` field that points to a different heap page. These are always non-HOT updates. [Table G.46](#) lists columns of this report table.

Table G.46. Top tables by new-page updated tuples

Column	Description
DB	Database name for the table
Tablespace	Name of the tablespace where the table is located
Schema	Schema name for the table
Table	Table name
NP Upd	Number of rows updated to a new heap page
%Upd	Number of new-page updated rows as the percentage of all rows updated
Upd	Number of rows updated, including HOT
Upd (HOT)	Number of rows HOT updated (i.e., with no separate index update required)

The report table “Top growing tables” shows top tables by growth. The data is based on the `pg_stat_all_tables` view. When there are no relation sizes collected with `pg_relation_size()`, relation-size estimates are based on the `pg_class.relpages` field. Since such values are less accurate, they are shown in square brackets. [Table G.47](#) lists columns of this report table.

Table G.47. Top growing tables

Column	Description	Field/Calculation
DB	Database name for the table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Size	Table size at the time of the last sample in the report interval	<code>pg_table_size()</code> - <code>pg_relation_size(toast)</code>
Growth	Table growth	
Ins	Number of rows inserted	<code>n_tup_ins</code>
Upd	Number of rows updated, including HOT	<code>n_tup_upd</code>
Del	Number of rows deleted	<code>n_tup_del</code>
Upd (HOT)	Number of rows HOT updated	<code>n_tup_hot_upd</code>

In the report table “Top indexes by blocks fetched”, blocks fetched include index blocks processed from disk (read) and from shared buffers (hit). The data is based on the `pg_statio_all_indexes` view. [Table G.48](#) lists columns of this report table.

Table G.48. Top indexes by blocks fetched

Column	Description	Field/Calculation
DB	Database name for the index	
Tablespace	Name of the tablespace where the index is located	
Schema	Schema name for the underlying table	
Table	Underlying table name	
Index	Index name	
Scans	Number of index scans initiated on this index	<code>idx_scan</code>
Blks	Number of blocks fetched for this index	<code>idx_blks_read</code> + <code>idx_blks_hit</code>
%Total	Percentage of blocks fetched for this index in all blocks fetched in the cluster	

The report table “Top indexes by blocks read” is also based on the `pg_statio_all_indexes` and `pg_stat_all_indexes` views. [Table G.49](#) lists columns of this report table.

Table G.49. Top indexes by blocks read

Column	Description	Field/Calculation
DB	Database name for the index	
Tablespace	Name of the tablespace where the index is located	
Schema	Schema name for the underlying table	
Table	Underlying table name	
Index	Index name	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Scans	Number of index scans initiated on this index	idx_scan
Blk Reads	Number of disk blocks read from this index	idx_blks_read
%Total	Percentage of disk blocks read from this index in all disk blocks read in the cluster	
Hits (%)	Percentage of index blocks got from buffers in all index blocks got for this index	

The report table “Top growing indexes” shows top indexes by growth. The table uses data from the [pg_stat_all_tables](#) and [pg_stat_all_indexes](#) views. When there are no relation sizes collected with [pg_relation_size\(\)](#), relation-size estimates are based on the [pg_class.relpages](#) field. Since such values are less accurate, they are shown in square brackets. [Table G.50](#) lists columns of this report table.

Table G.50. Top growing indexes

Column	Description	Field/Calculation
DB	Database name for the index	
Tablespace	Name of the tablespace where the index is located	
Schema	Schema name for the underlying table	
Table	Underlying table name	
Index	Index name	
Index Size	Index size at the time of the last sample in the report interval	pg_relation_size()
Index Growth	Index growth during the report interval	
Table Ins	Number of rows inserted into the underlying table	n_tup_ins
Table Upd	Number of rows updated in the underlying table, without HOT	n_tup_upd - n_tup_hot_upd
Table Del	Number of rows deleted from the underlying table	n_tup_del

The report table “Unused indexes” shows top non-scanned indexes (during the report interval) by DML operations on underlying tables that caused index support. Constraint indexes are not counted. The table uses data from the [pg_stat_all_tables](#) view. [Table G.51](#) lists columns of this report table.

Table G.51. Unused indexes

Column	Description	Field/Calculation
DB	Database name for the index	
Tablespace	Name of the tablespace where the index is located	
Schema	Schema name for the underlying table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Table	Underlying table name	
Index	Index name	
Index Size	Index size at the time of the last sample in the report interval	pg_relation_size()
Index Growth	Index growth during the report interval	
Table Ins	Number of rows inserted into the underlying table	n_tup_ins
Table Upd	Number of rows updated in the underlying table, without HOT	n_tup_upd - n_tup_hot_upd
Table Del	Number of rows deleted from the underlying table	n_tup_del

G.1.11.8. User function statistics

Tables in this section of the report show top functions in the cluster by statistics from the `pg_stat_user_functions` view. Times in the tables are provided in seconds.

The report table “Top functions by total time” shows top functions by the total time elapsed. The report table “Top functions by executions” shows top functions by the number of executions. The report table “Top trigger functions by total time” shows top trigger functions by the total time elapsed. [Table G.52](#) lists columns of these report tables.

Table G.52. User function statistics

Column	Description	Field/Calculation
DB	Database name for the function	
Schema	Schema name for the function	
Function	Function name	
Executions	Number of times this function has been called	calls
Total Time (s)	Total time spent in this function and all other functions called by it	total_time
Self Time (s)	Total time spent in this function itself, not including other functions called by it	self_time
Mean Time (s)	Mean time of a single function execution	total_time /calls
Mean self Time (s)	Mean self time of a single function execution	self_time /calls

G.1.11.9. Vacuum-related statistics

The report table “Top tables by vacuum time spent” is available if `pgpro_stats` can provide extended vacuum statistics. This table shows top tables by total time spent vacuuming them. The data is based on the `pgpro_stats_vacuum_tables` view. [Table G.53](#) lists columns of this report table.

Table G.53. Top tables by vacuum time spent

Column	Description	Field/Calculation
DB	Database name for the table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Total vacuum time	Total time of vacuuming this table	total_time
Delay vacuum time	Time spent sleeping in a vacuum delay point by vacuum operations performed on this table	delay_time
Read I/O time	Time spent reading database blocks by vacuum operations performed on this table	blk_read_time
Write I/O time	Time spent writing database blocks by vacuum operations performed on this table	blk_write_time
User CPU time	User CPU time of vacuuming tables of this database	user_time
System CPU time	System CPU time of vacuuming tables of this database	system_time
Vacuum count	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	vacuum_count
Autovacuum count	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count
Total fetched	Total number of database blocks fetched by vacuum operations performed on this table	total_blks_read + total_blks_hit
Heap fetched	Total number of blocks fetched from this table by vacuum operations performed on it	rel_blks_read + rel_blks_hit
Scanned	Number of pages examined by vacuum operations performed on this table	pages_scanned

The report table “Top indexes by vacuum time spent” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top indexes by total time spent vacuuming them. The data is based on the `pgpro_stats_vacuum_indexes` view. [Table G.54](#) lists columns of this report table.

Table G.54. Top indexes by vacuum time spent

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Index	Index name	
Total vacuum time	Total time of vacuuming this index	total_time

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Delay vacuum time	Time spent sleeping in a vacuum delay point by vacuum operations performed on this index	delay_time
Read I/O time	Time spent reading database blocks by vacuum operations performed on this index	blk_read_time
Write I/O time	Time spent writing database blocks by vacuum operations performed on this index	blk_write_time
User CPU time	User CPU time of vacuuming this index	user_time
System CPU time	System CPU time of vacuuming this index	system_time
Vacuum count	Number of times the underlying table has been manually vacuumed (not counting VACUUM FULL)	vacuum_count
Autovacuum count	Number of times the underlying table has been vacuumed by the autovacuum daemon	autovacuum_count
Total fetched	Total number of database blocks fetched by vacuum operations performed on the underlying table	total_blks_read + total_blks_hit
Index fetched	Total number of blocks fetched from this index by vacuum operations performed on it	rel_blks_read + rel_blks_hit

The report table “Top tables by blocks vacuum fetched” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top tables by blocks fetched vacuuming these tables. The data is based on the `pgpro_stats_vacuum_tables` view. [Table G.55](#) lists columns of this report table.

Table G.55. Top tables by blocks vacuum fetched

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
DB fetched	Total number of database blocks fetched by vacuum operations performed on this table	total_blks_read + total_blks_hit
%Total	Total number of database blocks fetched by vacuum operations performed on this table as the percentage of all blocks fetched in the cluster	
DB read	Total number of database blocks read by vacuum operations performed on this table	total_blks_read

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
%Total	Total number of database blocks read by vacuum operations performed on this table as the percentage of all blocks read in the cluster	
Heap fetched	Total number of blocks fetched from this table by vacuum operations performed on it	rel_blks_read + rel_blks_hit
%Rel	Total number of table blocks fetched by vacuum operations performed on this table as the percentage of all blocks fetched from this table	
Heap read	Total number of blocks read from this table by vacuum operations performed on it	rel_blks_read
%Rel	Total number of table blocks read by vacuum operations performed on this table as the percentage of all blocks read from this table	
Scanned	Number of pages examined by vacuum operations performed on this table	pages_scanned

The report table “Top indexes by blocks vacuum fetched” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top indexes by blocks fetched vacuuming underlying tables. The data is based on the `pgpro_stats_vacuum_indexes` view. [Table G.56](#) lists columns of this report table.

Table G.56. Top indexes by blocks vacuum fetched

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Index	Index name	
DB fetched	Total number of database blocks fetched by vacuum operations performed on this index	total_blks_read + total_blks_hit
%Total	Total number of database blocks fetched by vacuum operations performed on this index as the percentage of all blocks fetched in the cluster	
DB read	Total number of database blocks read by vacuum operations performed on this index	total_blks_read
%Total	Total number of database blocks read by vacuum operations per-	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	formed on this index as the percentage of all blocks read in the cluster	
Idx fetched	Total number of blocks fetched from this index by vacuum operations on it	rel_blks_read + rel_blks_hit
%Idx	Total number of index blocks fetched by vacuum operations performed on this index as the percentage of all blocks fetched from this index	
Idx read	Total number of blocks read from this index by vacuum operations performed on it	rel_blks_read
%Idx	Total number of index blocks read by vacuum operations performed on this index as the percentage of all blocks read from this index	

The report table “Top tables by blocks vacuum read” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top tables by blocks read vacuuming these tables. The data is based on the `pgpro_stats_vacuum_indexes` view. [Table G.57](#) lists columns of this report table.

Table G.57. Top tables by blocks vacuum read

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
DB read	Total number of database blocks read by vacuum operations performed on this table	total_blks_read
%Total	Total number of database blocks read by vacuum operations performed on this table as the percentage of all blocks read in the cluster	
%Hit	Total number of database blocks found in shared buffers by vacuum operations performed on this table as the percentage of all blocks fetched by vacuum operations performed on this table	
Heap read	Total number of database blocks vacuum operations read from this table	rel_blks_read
%Rel	Total number of table blocks read by vacuum operations performed on this table as the per-	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	centage of all blocks read from this table	
%Hit	Total number of table blocks found in shared buffers by vacuum operations performed on this table as the percentage of table blocks fetched by vacuum operations performed on this table	
Scanned	Number of pages examined by vacuum operations performed on this table	pages_scanned

The report table “Top indexes by blocks vacuum read” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top indexes by blocks read vacuuming underlying tables. The data is based on the `pgpro_stats_vacuum_indexes` view. [Table G.58](#) lists columns of this report table.

Table G.58. Top indexes by blocks vacuum read

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Index	Index name	
DB read	Total number of database blocks read by vacuum operations performed on this index	total_blks_read
%Total	Total number of database blocks read by vacuum operations performed on this index as the percentage of all blocks read in the cluster	
%Hit	Total number of database blocks found in shared buffers by vacuum operations performed on this index as the percentage of all database blocks fetched by vacuum operations performed on this index	
Idx read	Total number of blocks read from this index by vacuum operations performed on it	rel_blks_read
%Idx	Total number of index blocks read by vacuum operations performed on this index as the percentage of all blocks read from this index	
%Hit	Total number of index blocks found in shared buffers by vacuum operations performed on this index as the percentage of index	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
	blocks fetched by vacuum operations performed on this index	

The report table “Top tables by dead tuples vacuum left” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top tables by the number of dead tuples left by vacuum due to their visibility in transactions. The data is based on the `pgpro_stats_vacuum_tables` view. [Table G.59](#) lists columns of this report table.

Table G.59. Top tables by dead tuples vacuum left

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Dead tuples left	Total number of dead tuples vacuum operations left in this table due to their visibility in transactions	<code>dead_tuples</code>
Dead tuples deleted	Total number of dead tuples vacuum operations deleted from this table	<code>tuples_deleted</code>
%Eff	Vacuum efficiency in terms of deleted tuples. This is the percentage of dead tuples deleted from this table in all dead tuples to be deleted from this table.	$\text{tuples_deleted} * 100 / (\text{tuples_deleted} + \text{dead_tuples})$
Tuples del	Number of rows deleted	<code>pg_stat_all_tables .n_tup_del</code>
Tuples upd	Number of rows updated (includes HOT updated rows)	<code>pg_stat_all_tables .n_tup_upd</code>
Vacuum	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	<code>vacuum_count</code>
Autovacuum	Number of times this table has been vacuumed by the autovacuum daemon	<code>autovacuum_count</code>

The report table “Top tables by WAL size generated by vacuum” is available if [pgpro_stats](#) can provide extended vacuum statistics. This table shows top tables by the amount of WAL generated by vacuum operations performed on them. The data is based on the `pgpro_stats_vacuum_tables` view. [Table G.60](#) lists columns of this report table.

Table G.60. Top tables by WAL size generated by vacuum

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Table	Table name	
WAL size	Total amount of WAL bytes generated by vacuum operations performed on this table	wal_bytes
%Total	Total amount of WAL bytes generated by vacuum operations performed on this table as the percentage of all WAL generated in the cluster	
WAL FPI	Total number of WAL full page images generated by vacuum operations performed on this table	wal_fpi
Scanned blocks	Number of pages examined by vacuum operations performed on this table	pages_scanned
Dirtied blocks	Number of database blocks dirtied by vacuum operations performed on this table	total_blks_dirtied
Removed blocks	Number of pages removed by vacuum operations performed on this table	pages_removed
Vacuum	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	vacuum_count
Autovacuum	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count

The report table “Top tables by vacuum operations” shows top tables by the number of vacuum operations performed (`vacuum_count + autovacuum_count`). The data is based on the [pg_stat_all_tables](#) view. [Table G.61](#) lists columns of this report table.

Table G.61. Top tables by vacuum operations

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Vacuum count	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	vacuum_count
Autovacuum count	Number of times this table has been vacuumed by the autovacuum daemon	autovacuum_count
Ins	Number of rows inserted	n_tup_ins
Upd	Number of rows updated (includes HOT updated rows)	n_tup_upd
Del	Number of rows deleted	n_tup_del

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description	Field/Calculation
Upd (HOT)	Number of rows HOT updated	n_tup_hot_upd

The report table “Top tables by analyze operations” shows top tables by the number of analyze operations performed (`analyze_count + autoanalyze_count`). The data is based on the `pg_stat_all_tables` view. Table G.62 lists columns of this report table.

Table G.62. Top tables by analyze operations

Column	Description	Field/Calculation
DB	Database name for the table	
Tablespace	Name of the tablespace where the table is located	
Schema	Schema name for the table	
Table	Table name	
Analyze count	Number of times this table has been manually analyzed	<code>analyze_count</code>
Autoanalyze count	Number of times this table has been analyzed by the autovacuum daemon	<code>autoanalyze_count</code>
Ins	Number of rows inserted	<code>n_tup_ins</code>
Upd	Number of rows updated, including HOT	<code>n_tup_upd</code>
Del	Number of rows deleted	<code>n_tup_del</code>
Upd (HOT)	Number of rows HOT updated	<code>n_tup_hot_upd</code>

The report table “Top indexes by estimated vacuum load” shows top indexes by estimated implicit vacuum load. This load is calculated as the number of vacuum operations performed on the underlying table multiplied by the index size. The data is based on the `pg_stat_all_indexes` view. When there are no relation sizes collected with `pg_relation_size()`, relation-size estimates are based on the `pg_class.relpages` field. Since such values are less accurate, they are shown in square brackets. Table G.63 lists columns of this report table.

Table G.63. Top indexes by estimated vacuum load

Column	Description	Field/Calculation
DB	Database name for the index	
Tablespace	Name of the tablespace where the index is located	
Schema	Schema name for the underlying table	
Table	Underlying table name	
Index	Index name	
~Vacuum bytes	Vacuum load estimation	$(\text{vacuum_count} + \text{autovacuum_count}) * \text{index_size}$
Vacuum count	Number of times this table has been manually vacuumed (not counting <code>VACUUM FULL</code>)	<code>vacuum_count</code>
Autovacuum count	Number of times this table has been vacuumed by the autovacuum daemon	<code>autovacuum_count</code>

Column	Description	Field/Calculation
IX size	Average index size during the report interval	
Relsize	Average relation size during the report interval	

The report table “Top tables by dead tuples ratio” shows top tables larger than 5 MB by the ratio of dead tuples. Statistics are valid for the last sample in the report interval. The data is based on the [pg_stat_all_tables](#) view. When there are no relation sizes collected with [pg_relation_size\(\)](#), relation-size estimates are based on the [pg_class.relpages](#) field. Since such values are less accurate, they are shown in square brackets. [Table G.64](#) lists columns of this report table.

Table G.64. Top tables by dead tuples ratio

Column	Description	Field/Calculation
DB	Database name for the table	
Schema	Schema name for the table	
Table	Table name	
Live	Estimated number of live rows	n_live_tup
Dead	Estimated number of dead rows	n_dead_tup
%Dead	Percentage of dead rows in all table rows	
Last AV	Last time at which this table was vacuumed by the autovacuum daemon	last_autovacuum
Size	Table size	pg_table_size() - pg_relation_size(toast)

The report table “Top tables by modified tuples ratio” shows top tables larger than 5 MB by the ratio of modified tuples. Statistics are valid for the last sample in the report interval. The data is based on the [pg_stat_all_tables](#) view. When there are no relation sizes collected with [pg_relation_size\(\)](#), relation-size estimates are based on the [pg_class.relpages](#) field. Since such values are less accurate, they are shown in square brackets. [Table G.65](#) lists columns of this report table.

Table G.65. Top tables by modified tuples ratio

Column	Description	Field/Calculation
DB	Database name for the table	
Schema	Schema name for the table	
Table	Table name	
Live	Estimated number of live rows	n_live_tup
Dead	Estimated number of dead rows	n_dead_tup
Mod	Estimated number of rows modified since this table was last analyzed	n_mod_since_analyze
%Mod	Percentage of modified rows in all table rows	
Last AA	Last time at which this table was analyzed by the autovacuum daemon	last_autoanalyze
Size	Table size	pg_table_size() - pg_relation_size(toast)

G.1.11.10. Cluster settings during the report interval

This section of the report contains a table with Postgres Pro GUC parameters, values of functions `version()`, `pg_postmaster_start_time()`, `pg_conf_load_time()` and the `system_identifier` field of the `pg_control_system()` function during the report interval. The data in the table is grouped under `Defined settings` and `Default settings`. [Table G.66](#) lists columns of this report table.

Table G.66. Cluster settings during the report interval

Column	Description
Setting	Name of the parameter
reset_val	<code>reset_val</code> field of the <code>pg_settings</code> view. Settings changed during the report interval are shown in bold font.
Unit	Unit of the setting
Source	Configuration file where this setting is defined, semicolon, line number
Notes	Timestamp of the sample where this value was first observed

G.1.11.11. Extension versions during the report interval

This section of the report contains a table that lists installed extension versions found in databases during the report interval. `First seen` and `Last seen` columns are not shown if the extension versions have not changed during the report interval. [Table G.67](#) lists columns of this report table.

Table G.67. Extension versions during the report interval

Column	Description
Name	Extension name
DB	Database name
First seen	Timestamp of the sample where this extension version appeared first
Last seen	Timestamp of the sample where this extension version appeared last
Version	Version name of the extension

G.1.12. pgpro_pwr Diagnostic Tools

`pgpro_pwr` provides self-diagnostic tools.

G.1.12.1. Collecting Detailed Timing Statistics for Sampling Procedures

`pgpro_pwr` collects detailed timing statistics of taking samples when the `pgpro_pwr.track_sample_timings` parameter is on. You can get the results from the `v_sample_timings` view. [Table G.68](#) lists columns of this view.

Table G.68. `v_sample_timings` View

Column	Description
<code>server_name</code>	Name of the server
<code>sample_id</code>	Sample identifier
<code>sample_time</code>	Time when the sample was taken
<code>sampling_event</code>	Sampling stage. See Table G.69 for descriptions of sampling stages.

Postgres Pro Modules
and Extensions Shipped
as Individual Packages

Column	Description
time_spent	Time spent in the event

Table G.69. `sampling_event` Description

Event	Description
total	Taking the sample (all stages)
connect	Making dblink connection to the server
get server environment	Getting server GUC parameters, available extensions, etc.
collect database stats	Querying the <code>pg_stat_database</code> view for statistics on databases
calculate database stats	Calculating differential statistics on databases since the previous sample
collect tablespace stats	Querying the <code>pg_tablespace</code> view for statistics on tablespaces
collect statement stats	Collecting statistics on statements using the <code>pg-pro_stats</code> and <code>pg_stat_kcache</code> extensions
query pg_stat_bgwriter	Collecting cluster statistics using the <code>pg_stat_bgwriter</code> view
query pg_stat_archiver	Collecting cluster statistics using the <code>pg_stat_archiver</code> view
collect object stats	Collecting statistics on database objects. Includes events from Table G.70 .
maintain repository	Executing support routines
calculate tablespace stats	Calculating differential statistics on tablespaces
calculate object stats	Calculating differential statistics on database objects. Includes events from Table G.71 .
calculate cluster stats	Calculating cluster differential statistics
calculate archiver stats	Calculating archiver differential statistics
delete obsolete samples	Deleting obsolete baselines and samples

Table G.70. Events of Collecting Statistics on Database Objects

Event	Description
<code>db:dbname collect tables stats</code>	Collecting statistics on tables for the <code>dbname</code> database
<code>db:dbname collect indexes stats</code>	Collecting statistics on indexes for the <code>dbname</code> database
<code>db:dbname collect functions stats</code>	Collecting statistics on functions for the <code>dbname</code> database

Table G.71. Events of Calculating Differences of Statistics on Database Objects

Event	Description
<code>calculate tables stats</code>	Calculating differential statistics on tables of all databases
<code>calculate indexes stats</code>	Calculating differential statistics on indexes of all databases

Event	Description
calculate functions stats	Calculating differential statistics on functions of all databases

G.1.13. Important Notes

When using the pgpro_pwr extension, be aware of the following:

- Postgres Pro collects execution statistics *after* the execution is complete. If a single execution of a statement lasts for several samples, it will only affect statistics of the last sample (in which the execution completed). Besides, statistics on statements that are still running are unavailable. Maintenance processes, such as vacuum and checkpointing, will update the statistics only on completion.
- Resetting any Postgres Pro statistics may affect the accuracy of the next sample.
- Exclusive locks on relations conflict with calculation of the relation size. If the `take_sample()` function is unable to acquire a lock for a short period of time (3 seconds), it will fail and no sample will be generated.

Appendix H. Additional Supplied Programs

This appendix, [Appendix F](#), and [Appendix G](#) contain information regarding additional modules available in the Postgres Pro Shardman distribution. See [Appendix F](#) and [Appendix G](#) for more information about the server extensions and plug-ins.

This appendix covers the utility programs. Once installed, they are found in the `bin` directory of the Postgres Pro Shardman installation and can be used like any other program.

H.1. Additional PostgreSQL/Postgres Pro Client Applications

This section covers PostgreSQL/Postgres Pro additional client applications included in the Postgres Pro Shardman distribution. They can be run from anywhere, independent of where the database server resides. See also [Postgres Pro Client Applications](#) for information about client applications that are part of the core Postgres Pro Shardman distribution. Note also the [ora2pgpro](#) utility, which is developed independently from the core project.

oid2name

oid2name — resolve OIDs and file nodes in a Postgres Pro data directory

Synopsis

```
oid2name [option...]
```

Description

oid2name is a utility program that helps administrators to examine the file structure used by Postgres Pro. To make use of it, you need to be familiar with the database file structure, which is described in [Chapter 66](#).

Note

The name “oid2name” is historical, and is actually rather misleading, since most of the time when you use it, you will really be concerned with tables' filenode numbers (which are the file names visible in the database directories). Be sure you understand the difference between table OIDs and table filenodes!

oid2name connects to a target database and extracts OID, filenode, and/or table name information. You can also have it show database OIDs or tablespace OIDs.

Options

oid2name accepts the following command-line arguments:

```
-f filenode
```

```
--filenode=filenode
```

show info for table with filenode *filenode*.

```
-i
```

```
--indexes
```

include indexes and sequences in the listing.

```
-o oid
```

```
--oid=oid
```

show info for table with OID *oid*.

```
-q
```

```
--quiet
```

omit headers (useful for scripting).

```
-s
```

```
--tablespaces
```

show tablespace OIDs.

```
-S
```

```
--system-objects
```

include system objects (those in `information_schema`, `pg_toast` and `pg_catalog` schemas).

```
-t tablename_pattern
```

```
--table=tablename_pattern
```

show info for table(s) matching *tablename_pattern*.

`-V`
`--version`
Print the oid2name version and exit.

`-x`
`--extended`
display more information about each object shown: tablespace name, schema name, and OID.

`-?`
`--help`
Show help about oid2name command line arguments, and exit.

oid2name also accepts the following command-line arguments for connection parameters:

`-d database`
`--dbname=database`
database to connect to.

`-h host`
`--host=host`
database server's host.

`-H host`
database server's host. Use of this parameter is *deprecated* as of Postgres Pro 12.

`-p port`
`--port=port`
database server's port.

`-U username`
`--username=username`
user name to connect as.

To display specific tables, select which tables to show by using `-o`, `-f` and/or `-t`. `-o` takes an OID, `-f` takes a filenode, and `-t` takes a table name (actually, it's a `LIKE` pattern, so you can use things like `foo%`). You can use as many of these options as you like, and the listing will include all objects matched by any of the options. But note that these options can only show objects in the database given by `-d`.

If you don't give any of `-o`, `-f` or `-t`, but do give `-d`, it will list all tables in the database named by `-d`. In this mode, the `-s` and `-i` options control what gets listed.

If you don't give `-d` either, it will show a listing of database OIDs. Alternatively you can give `-s` to get a tablespace listing.

Environment

`PGHOST`
`PGPORT`
`PGUSER`

Default connection parameters.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by libpq (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

Notes

oid2name requires a running database server with non-corrupt system catalogs. It is therefore of only limited use for recovering from catastrophic database corruption situations.

Examples

```
$ # what's in this database server, anyway?
$ oid2name
All databases:
  Oid  Database Name  Tablespace
-----
 17228      alvherre  pg_default
 17255      regression pg_default
 17227      template0  pg_default
    1      template1  pg_default

$ oid2name -s
All tablespaces:
  Oid  Tablespace Name
-----
 1663      pg_default
 1664      pg_global
155151      fastdisk
155152      bigdisk

$ # OK, let's look into database alvherre
$ cd $PGDATA/base/17228

$ # get top 10 db objects in the default tablespace, ordered by size
$ ls -lS * | head -10
-rw----- 1 alvherre alvherre 136536064 sep 14 09:51 155173
-rw----- 1 alvherre alvherre 17965056 sep 14 09:51 1155291
-rw----- 1 alvherre alvherre 1204224 sep 14 09:51 16717
-rw----- 1 alvherre alvherre 581632 sep 6 17:51 1255
-rw----- 1 alvherre alvherre 237568 sep 14 09:50 16674
-rw----- 1 alvherre alvherre 212992 sep 14 09:51 1249
-rw----- 1 alvherre alvherre 204800 sep 14 09:51 16684
-rw----- 1 alvherre alvherre 196608 sep 14 09:50 16700
-rw----- 1 alvherre alvherre 163840 sep 14 09:50 16699
-rw----- 1 alvherre alvherre 122880 sep 6 17:51 16751

$ # What file is 155173?
$ oid2name -d alvherre -f 155173
From database "alvherre":
  Filenode  Table Name
-----
 155173      accounts

$ # you can ask for more than one object
$ oid2name -d alvherre -f 155173 -f 1155291
From database "alvherre":
  Filenode  Table Name
-----
 155173      accounts
1155291  accounts_pkey

$ # you can mix the options, and get more details with -x
$ oid2name -d alvherre -t accounts -f 1155291 -x
```

From database "alvherre":

Filenode	Table Name	Oid	Schema	Tablespace
155173	accounts	155173	public	pg_default
1155291	accounts_pkey	1155291	public	pg_default

\$ # show disk space for every db object

\$ du [0-9]* |

> while read SIZE FILENODE

> do

> echo "\$SIZE `oid2name -q -d alvherre -i -f \$FILENODE`"

> done

16 1155287 branches_pkey

16 1155289 tellers_pkey

17561 1155291 accounts_pkey

...

\$ # same, but sort by size

\$ du [0-9]* | sort -rn | while read SIZE FN

> do

> echo "\$SIZE `oid2name -q -d alvherre -f \$FN`"

> done

133466 155173 accounts

17561 1155291 accounts_pkey

1177 16717 pg_proc_proname_args_nsp_index

...

\$ # If you want to see what's in tablespaces, use the pg_tblspc directory

\$ cd \$PGDATA/pg_tblspc

\$ oid2name -s

All tablespaces:

Oid	Tablespace Name
1663	pg_default
1664	pg_global
155151	fastdisk
155152	bigdisk

\$ # what databases have objects in tablespace "fastdisk"?

\$ ls -d 155151/*

155151/17228/ 155151/PG_VERSION

\$ # Oh, what was database 17228 again?

\$ oid2name

All databases:

Oid	Database Name	Tablespace
17228	alvherre	pg_default
17255	regression	pg_default
17227	template0	pg_default
1	template1	pg_default

\$ # Let's see what objects does this database have in the tablespace.

\$ cd 155151/17228

\$ ls -l

total 0

-rw----- 1 postgres postgres 0 sep 13 23:20 155156

```
$ # OK, this is a pretty small table ... but which one is it?
$ oid2name -d alvherre -f 155156
From database "alvherre":
  Filenode  Table Name
-----
    155156         foo
```

Author

B. Palmer <bpalmer@crimelabs.net>

pg_probackup

pg_probackup — manage backup and recovery of Postgres Pro Shardman database clusters

Synopsis

```
pg_probackup version

pg_probackup help [command]

pg_probackup init -B backup_dir --skip-if-exists

pg_probackup add-instance -B backup_dir -D data_dir --instance instance_name --skip-if-exists

pg_probackup del-instance -B backup_dir --instance instance_name

pg_probackup set-config -B backup_dir --instance instance_name [option...]

pg_probackup set-backup -B backup_dir --instance instance_name -i backup_id [option...]

pg_probackup show-config -B backup_dir --instance instance_name [option...]

pg_probackup show -B backup_dir [option...]

pg_probackup backup -B backup_dir --instance instance_name -b backup_mode [option...]

pg_probackup restore -B backup_dir --instance instance_name [option...]

pg_probackup checkdb -B backup_dir --instance instance_name -D data_dir [option...]

pg_probackup validate -B backup_dir [option...]

pg_probackup merge -B backup_dir --instance instance_name -i backup_id [option...]

pg_probackup delete -B backup_dir --instance instance_name { -i backup_id | --delete-wal | --delete-expired | --merge-expired } [option...]

pg_probackup archive-push -B backup_dir --instance instance_name --wal-file-path wal_file_path --wal-file-name wal_file_name [option...]

pg_probackup archive-get -B backup_dir --instance instance_name --wal-file-path wal_file_path --wal-file-name wal_file_name [option...]

pg_probackup catchup -b catchup_mode --source-pgdata=path_to_pgdata_on_remote_server --destination-pgdata=path_to_local_dir [option...]
```

Description

pg_probackup is a utility to manage backup and recovery of Postgres Pro database clusters. It is designed to perform periodic backups of the Postgres Pro instance that enable you to restore the server in case of a failure. pg_probackup supports PostgreSQL 11 or higher. In Postgres Pro Shardman, pg_probackup provides S3 (Simple Storage Service) support for storing data in private clouds.

Note

pg_probackup provides complete processing of S3 interface logging.

- [Overview](#)
- [Quick Start](#)
- [Installation and Setup](#)
- [Command-Line Reference](#)
- [Usage](#)

Overview

As compared to other backup solutions, `pg_probackup` offers the following benefits that can help you implement different backup strategies and deal with large amounts of data:

- S3 support for storing data in private clouds using MinIO object storage, Amazon S3 storage, and VK Cloud storage: provided in Postgres Pro Shardman. Backup data is transferred to and from S3 without saving it in intermediate locations thus eliminating the need of having a large temporary storage.
- Incremental backup: with three different incremental modes, you can plan the backup strategy in accordance with your data flow. Incremental backups allow you to save disk space and speed up backup as compared to taking full backups. It is also faster to restore the cluster by applying incremental backups than by replaying WAL files.
- Incremental restore: speed up restore from a backup by reusing valid unchanged pages available in PGDATA.
- CFS (Compressed File System) support for incremental backups in `DELTA`, `PAGE`, and `PTRACK` (the fastest) modes: provided in Postgres Pro Shardman.
- Validation: automatic data consistency checks and on-demand backup validation without actual data recovery.
- Verification: on-demand verification of Postgres Pro instance with the `checkdb` command.
- Retention: managing WAL archive and backups in accordance with retention policy. You can configure retention policy based on recovery time or the number of backups to keep, as well as specify time to live (TTL) for a particular backup. Expired backups can be merged or deleted.
- Parallelization: running `backup`, `restore`, `merge`, `delete`, `validate`, and `checkdb` processes on multiple parallel threads.
- Compression: storing backup data in a compressed state to save disk space.
- Deduplication: saving disk space by excluding non-data files (such as `_vm` or `_fsm`) from incremental backups if these files have not changed since they were copied into one of the previous backups in this incremental chain.
- Remote operations: backing up Postgres Pro instance located on a remote system or restoring a backup remotely.
- Backups from a standby: avoiding extra load on the primary by taking backups from a standby server.
- External directories: backing up files and directories located outside of the Postgres Pro data directory (PGDATA), such as scripts, configuration files, logs, or SQL dump files.
- Backup catalog: getting the list of backups and the corresponding meta information in plain text or JSON formats.
- Archive catalog: getting the list of all WAL timelines and the corresponding meta information in plain text or JSON formats.
- Partial restore: restoring only the specified databases.
- Catchup: cloning a Postgres Pro instance for a fallen-behind standby server to “catch up” with the primary.

To manage backup data, `pg_probackup` creates a *backup catalog*. This is a directory that stores all backup files with additional meta information, as well as WAL archives required for point-in-time recovery. You can store backups for different instances in separate subdirectories of a single backup catalog.

Using `pg_probackup`, you can take full or incremental [backups](#):

- FULL backups contain all the data files required to restore the database cluster.
- Incremental backups operate at the page level, only storing the data that has changed since the previous backup. It allows you to save disk space and speed up the backup process as compared to

taking full backups. It is also faster to restore the cluster by applying incremental backups than by replaying WAL files. `pg_probackup` supports the following modes of incremental backups:

- **DELTA backup.** In this mode, `pg_probackup` reads all data files in the data directory and copies only those pages that have changed since the previous backup. This mode can impose read-only I/O pressure equal to a full backup.
- **PAGE backup.** In this mode, `pg_probackup` scans all WAL files in the archive from the moment the previous full or incremental backup was taken. Newly created backups contain only the pages that were mentioned in WAL records. This requires all the WAL files since the previous backup to be present in the WAL archive. If the size of these files is comparable to the total size of the database cluster files, speedup is smaller, but the backup still takes less space. You have to configure WAL archiving as explained in [Setting up continuous WAL archiving](#) to make PAGE backups.
- **PTRACK backup.** In this mode, Postgres Pro tracks page changes on the fly. Continuous archiving is not necessary for it to operate. Each time a relation page is updated, this page is marked in a special PTRACK bitmap. Tracking implies some minor overhead on the database server operation, but speeds up incremental backups significantly.

`pg_probackup` can take only physical online backups, and online backups require WAL for consistent recovery. So regardless of the chosen backup mode (FULL, PAGE or DELTA), any backup taken with `pg_probackup` must use one of the following *WAL delivery modes*:

- **ARCHIVE.** Such backups rely on [continuous archiving](#) to ensure consistent recovery. This is the default WAL delivery mode.
- **STREAM.** Such backups include all the files required to restore the cluster to a consistent state at the time the backup was taken. Regardless of [continuous archiving](#) having been set up or not, the WAL segments required for consistent recovery are streamed via replication protocol during backup and included into the backup files. That's why such backups are called *autonomous*, or *stand-alone*.

Limitations

`pg_probackup` currently has the following limitations:

- The remote mode is not supported on Windows systems.
- On Unix systems, for Postgres Pro, a backup can be made only by the same OS user that has started the Postgres Pro server. For example, if Postgres Pro server is started by user `postgres`, the `backup` command must also be run by user `postgres`. To satisfy this requirement when taking backups in the [remote mode](#) using SSH, you must set `--remote-user` option to `postgres`.
- The Postgres Pro server from which the backup was taken and the restored server must be compatible by the [block_size](#) and [wal_block_size](#) parameters and have the same major release number. Depending on cluster configuration, Postgres Pro itself may apply additional restrictions, such as CPU architecture or `libc/icu` versions.
- Special limitations of `pg_probackup` in Postgres Pro Shardman:
 - Incremental backup and restore for CFS tablespaces require `ptrack` 2.4.0 or higher.
 - Only the following commands can be launched in the remote mode: [add-instance](#), [backup](#), [restore](#), [delete](#), [catchup](#), [archive-push](#), and [archive-get](#).

Quick Start

To quickly get started with `pg_probackup`, complete the steps below. This will set up FULL and DELTA backups in the remote mode and demonstrate some basic `pg_probackup` operations. In the following, these terms are used:

- `backup` — Postgres Pro role used to connect to the Postgres Pro cluster.
- `backupdb` — database used to connect to the Postgres Pro cluster.
- `backup_host` — host with the backup catalog.
- `backup_user` — user on `backup_host` running all `pg_probackup` operations.
- `/mnt/backups` — directory on `backup_host` where the backup catalog is stored.

- `postgres_host` — host with the Postgres Pro cluster.
- `postgres` — user on `postgres_host` under which Postgres Pro cluster processes are running.
- `/var/lib/pgpro/std-16/data` — Postgres Pro data directory on `postgres_host`.

Steps to perform:

1. Install `pg_probackup` on both `backup_host` and `postgres_host`.
2. [Set up an SSH connection](#) from `backup_host` to `postgres_host`.
3. [Configure](#) your database cluster for **STREAM backups**.
4. Initialize the backup catalog:

```
backup_user@backup_host:~$ pg_probackup init -B /mnt/backups
INFO: Backup catalog '/mnt/backups' successfully initialized
```

5. Add a backup instance called `mydb` to the backup catalog:

```
backup_user@backup_host:~$ pg_probackup add-instance \
-B /mnt/backups \
-D /var/lib/pgpro/std-16/data \
--instance=node \
--remote-host=postgres_host \
--remote-user=postgres
INFO: Instance 'node' successfully initialized
```

6. Make a FULL backup:

```
backup_user@backup_host:~$ pg_probackup backup \
-B /mnt/backups \
-b FULL \
--instance=node \
--stream \
--compress-algorithm=zstd \
--remote-host=postgres_host \
--remote-user=postgres \
-U backup \
-d backupdb
INFO: Backup start, pg_probackup version: 2.7.3, instance: node, backup ID: SBOL6J,
backup mode: FULL, wal mode: STREAM, remote: true, compress-algorithm: zstd,
compress-level: 1
WARNING: pgpro_edition() function is old-style and will be removed in future major
release, use pgpro_edition GUC variable instead.
INFO: This PostgreSQL instance was initialized with data block checksums. Data
block corruption will be detected
INFO: Database backup start
INFO: wait for pg_backup_start()
INFO: PGDATA size: 96MB
INFO: Current Start LSN: 0/8000028, TLI: 1
INFO: Start transferring data files
INFO: Data files are transferred, time elapsed: 2s
INFO: wait for pg_stop_backup()
INFO: pg_stop_backup() successfully executed
INFO: stop_stream_lsn 0/9000000 currentpos 0/9000000
INFO: backup->stop_lsn 0/8004E40
INFO: Getting the Recovery Time from WAL
INFO: Syncing backup files to disk
INFO: Backup files are synced, time elapsed: 1s
INFO: Validating backup SBOL6J
INFO: Backup SBOL6J data files are valid
INFO: Backup SBOL6J resident size: 53MB
```



```
INFO: Backup SBOL6J completed
```

7. List the backups of the instance:

```
backup_user@backup_host:~$ pg_probackup show \  
-B /mnt/backups \  
--instance=node
```

```
=====
```

Instance	Version	ID	Recovery Time	Mode	WAL Mode	TLI
Time	Data	WAL	Zalg Zratio	Start LSN	Stop LSN	Status
node	17	SBOL6J	2024-04-09 18:18:21.970314+03	FULL	STREAM	1/0
4s	37MB	16MB	zstd 2.57	0/8000028	0/8004E40	OK

```
=====
```

8. Make an incremental backup in the DELTA mode:

```
backup_user@backup_host:~$ pg_probackup backup \  
-B /mnt/backups \  
-b DELTA \  
--instance=node \  
--stream \  
--compress-algorithm=zstd \  
--remote-host=postgres_host \  
--remote-user=postgres \  
-U backup \  
-d backupdb
```

```
INFO: Backup start, pg_probackup version: 2.7.3, instance: node, backup ID: SBOL6N,  
backup mode: DELTA, wal mode: STREAM, remote: true, compress-algorithm: zstd,  
compress-level: 1
```

```
WARNING: pgpro_edition() function is old-style and will be removed in future major  
release, use pgpro_edition GUC variable instead.
```

```
INFO: This PostgreSQL instance was initialized with data block checksums. Data  
block corruption will be detected
```

```
INFO: Database backup start
```

```
INFO: wait for pg_backup_start()
```

```
INFO: Parent backup: SBOL6J
```

```
INFO: PGDATA size: 96MB
```

```
INFO: Current Start LSN: 0/9000028, TLI: 1
```

```
INFO: Parent Start LSN: 0/8000028, TLI: 1
```

```
INFO: Start transferring data files
```

```
INFO: Data files are transferred, time elapsed: 1s
```

```
INFO: wait for pg_stop_backup()
```

```
INFO: pg_stop_backup() successfully executed
```

```
INFO: stop_stream_lsn 0/A000000 currentpos 0/A000000
```

```
INFO: backup->stop_lsn 0/9000190
```

```
INFO: Getting the Recovery Time from WAL
```

```
INFO: Syncing backup files to disk
```

```
INFO: Backup files are synced, time elapsed: 0
```

```
INFO: Validating backup SBOL6N
```

```
INFO: Backup SBOL6N data files are valid
```

```
INFO: Backup SBOL6N resident size: 17MB
```

```
INFO: Backup SBOL6N completed
```

9. Add or modify some parameters in the pg_probackup configuration file, so that you do not have to specify them each time on the command line:

```
backup_user@backup_host:~$ pg_probackup set-config \  
-B /mnt/backups \  
--instance=node \  
--remote-host=postgres_host \  

```

```
--remote-user=postgres \  
-U backup \  
-d backupdb
```

10. Check the configuration of the instance:

```
backup_user@backup_host:~$ pg_probackup show-config \  
-B /mnt/backups \  
--instance=node  
# Backup instance information  
pgdata = /var/lib/pgpro/std-16/data  
system-identifier = 7355886958826772732  
xlog-seg-size = 16777216  
# Connection parameters  
pgdatabase = backupdb  
pghost = postgres_host  
pguser = backup  
# Archive parameters  
archive-timeout = 5min  
# Logging parameters  
log-level-console = INFO  
log-level-file = OFF  
log-format-console = PLAIN  
log-format-file = PLAIN  
log-filename = pg_probackup.log  
log-rotation-size = 0TB  
log-rotation-age = 0d  
# Retention parameters  
retention-redundancy = 0  
retention-window = 0  
wal-depth = 0  
# Compression parameters  
compress-algorithm = none  
compress-level = 1  
# Remote access parameters  
remote-proto = ssh  
remote-host = postgres_host  
remote-user = postgres
```

Note that the parameters not modified via `set-config` retain their default values.

11. Make another incremental backup in the DELTA mode, omitting the parameters stored in the configuration file earlier:

```
backup_user@backup_host:~$ pg_probackup backup \  
-B /mnt/backups \  
-b DELTA \  
--instance=node \  
--stream \  
--compress-algorithm=zstd  
INFO: Backup start, pg_probackup version: 2.7.3, instance: node, backup ID: SBOL6P,  
backup mode: DELTA, wal mode: STREAM, remote: true, compress-algorithm: zstd,  
compress-level: 1  
WARNING: pgpro_edition() function is old-style and will be removed in future major  
release, use pgpro_edition GUC variable instead.  
INFO: This PostgreSQL instance was initialized with data block checksums. Data  
block corruption will be detected  
INFO: Database backup start  
INFO: wait for pg_backup_start()  
INFO: Parent backup: SBOL6N
```

```

INFO: PGDATA size: 96MB
INFO: Current Start LSN: 0/A000028, TLI: 1
INFO: Parent Start LSN: 0/9000028, TLI: 1
INFO: Start transferring data files
INFO: Data files are transferred, time elapsed: 1s
INFO: wait for pg_stop_backup()
INFO: pg_stop_backup() successfully executed
INFO: stop_stream_lsn 0/B000000 currentpos 0/B000000
INFO: backup->stop_lsn 0/A000190
INFO: Getting the Recovery Time from WAL
INFO: Syncing backup files to disk
INFO: Backup files are synced, time elapsed: 0
INFO: Validating backup SBOL6P
INFO: Backup SBOL6P data files are valid
INFO: Backup SBOL6P resident size: 17MB
INFO: Backup SBOL6P completed

```

12. List the backups of the instance again:

```

backup_user@backup_host:~$ pg_probackup show \
    -B /mnt/backups \
    --instance=node

```

Instance	Version	ID	Recovery Time	Mode	WAL Mode	TLI	
Time	Data	WAL	Zalg	Zratio	Start LSN	Stop LSN	Status
node	17	SBOL6P	2024-04-09 18:18:26.630175+03	DELTA	STREAM	1/1	
1s	1147kB	16MB	zstd	1.00	0/A000028	0/A000190	OK
node	17	SBOL6N	2024-04-09 18:18:25.015713+03	DELTA	STREAM	1/1	
2s	1160kB	16MB	zstd	1.04	0/9000028	0/9000190	OK
node	17	SBOL6J	2024-04-09 18:18:21.970314+03	FULL	STREAM	1/0	
4s	37MB	16MB	zstd	2.57	0/8000028	0/8004E40	OK

13. Restore the data from the latest available backup to an arbitrary location:

```

backup_user@backup_host:~$ pg_probackup restore \
    -B /mnt/backups \
    -D /var/lib/pgpro/std-16/staging-data \
    --instance=node
INFO: Validating parents for backup SBOL6P
INFO: Validating backup SBOL6J
INFO: Backup SBOL6J data files are valid
INFO: Validating backup SBOL6N
INFO: Backup SBOL6N data files are valid
INFO: Validating backup SBOL6P
INFO: Backup SBOL6P data files are valid
INFO: Backup SBOL6P WAL segments are valid
INFO: Backup SBOL6P is valid.
INFO: Restoring the database from the backup starting at 2024-04-09 18:18:25+03 on
    localhost
INFO: Start restoring backup files. PGDATA size: 112MB
INFO: Backup files are restored. Transferred bytes: 112MB, time elapsed: 2s
INFO: Restore incremental ratio (less is better): 100% (112MB/112MB)
INFO: Syncing restored files to disk
INFO: Restored backup files are synced, time elapsed: 1s
INFO: Restore of backup SBOL6P completed.

```

Installation and Setup

Once you have pg_probackup installed, complete the following setup:

- Initialize the backup catalog.
- Add a new backup instance to the backup catalog.
- Configure the database cluster to enable `pg_probackup` backups.
- Optionally, configure SSH for running `pg_probackup` operations in the remote mode.
- Optionally, configure S3 for running `pg_probackup` connected to the S3 storage.

Initializing the Backup Catalog

`pg_probackup` stores all WAL and backup files in the corresponding subdirectories of the backup catalog.

To initialize the backup catalog, run the following command:

```
pg_probackup init -B backup_dir
```

where *backup_dir* is the path to the backup catalog. If the *backup_dir* already exists, it must be empty. Otherwise, `pg_probackup` returns an error.

The user launching `pg_probackup` must have full access to the *backup_dir* directory.

`pg_probackup` creates the *backup_dir* backup catalog, with the following subdirectories:

- `wal/` — directory for WAL files.
- `backups/` — directory for backup files.

Once the backup catalog is initialized, you can add a new backup instance.

Adding a New Backup Instance

`pg_probackup` can store backups for multiple database clusters in a single backup catalog. To set up the required subdirectories, you must add a backup instance to the backup catalog for each database cluster you are going to back up.

To add a new backup instance, run the following command:

```
pg_probackup add-instance -B backup_dir -D data_dir --instance=instance_name  
[remote_options]
```

Where:

- *data_dir* is the data directory of the cluster you are going to back up. To set up and use `pg_probackup`, write access to this directory is required.
- *instance_name* is the name of the subdirectories that will store WAL and backup files for this cluster.
- [remote_options](#) are optional parameters that need to be specified only if *data_dir* is located on a remote system.

`pg_probackup` creates the *instance_name* subdirectories under the `backups/` and `wal/` directories of the backup catalog. The `backups/instance_name` directory contains the `pg_probackup.conf` configuration file that controls `pg_probackup` settings for this backup instance. If you run this command with the [remote_options](#), the specified parameters will be added to `pg_probackup.conf`.

For details on how to fine-tune `pg_probackup` configuration, see [the section called “Configuring pg_probackup”](#).

The user launching `pg_probackup` must have full access to *backup_dir* directory and at least read-only access to *data_dir* directory. If you specify the path to the backup catalog in the `BACKUP_PATH` environment variable, you can omit the corresponding option when running `pg_probackup` commands.

Note

For Postgres Pro 11 or higher, it is recommended to use the [group access](#) feature, so that backups can be done by any OS user in the same group as the cluster owner. In this case, the user should have read permissions for the cluster directory.

Configuring the Database Cluster

Although `pg_probackup` can be used by a superuser, it is recommended to create a separate role with the minimum permissions required for the chosen backup strategy. In these configuration instructions, the `backup` role is used as an example.

For security reasons, it is recommended to run the configuration SQL queries below in a separate database.

```
postgres=# CREATE DATABASE backupdb;
postgres=# \c backupdb
```

To perform a [backup](#), the following permissions for role `backup` are required only in the database **used for connection** to the Postgres Pro server.

For Postgres Pro versions 11 — 14:

```
BEGIN;
CREATE ROLE backup WITH LOGIN;
GRANT USAGE ON SCHEMA pg_catalog TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.current_setting(text) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.set_config(text, text, boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_is_in_recovery() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_start_backup(text, boolean, boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_stop_backup(boolean, boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_create_restore_point(text) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_switch_wal() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_last_wal_replay_lsn() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_current() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_current_snapshot() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_snapshot_xmax(txid_snapshot) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_control_checkpoint() TO backup;
COMMIT;
```

For Postgres Pro 15 or higher:

```
BEGIN;
CREATE ROLE backup WITH LOGIN;
GRANT USAGE ON SCHEMA pg_catalog TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.current_setting(text) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.set_config(text, text, boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_is_in_recovery() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_backup_start(text, boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_backup_stop(boolean) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_create_restore_point(text) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_switch_wal() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_last_wal_replay_lsn() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_current() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_current_snapshot() TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.txid_snapshot_xmax(txid_snapshot) TO backup;
GRANT EXECUTE ON FUNCTION pg_catalog.pg_control_checkpoint() TO backup;
COMMIT;
```

In the `pg_hba.conf` file, allow connection to the database cluster on behalf of the `backup` role.

Since `pg_probackup` needs to read cluster files directly, `pg_probackup` must be started by (or connected to, if used in the remote mode) the OS user that has read access to all files and directories inside the data directory (`PGDATA`) you are going to back up.

Depending on whether you plan to take [standalone](#) or [archive](#) backups, Postgres Pro cluster configuration will differ, as specified in the sections below. To back up the database cluster from a standby server, run `pg_probackup` in the remote mode, or create PTRACK backups, additional setup is required.

For details, see the sections [Setting up STREAM Backups](#), [Setting up continuous WAL archiving](#), [Setting up Backup from Standby](#), [Configuring the Remote Mode](#), [Setting up Partial Restore](#), and [Setting up PTRACK Backups](#).

Setting up STREAM Backups

To set up the cluster for [STREAM](#) backups, complete the following steps:

- If the `backup` role does not exist, create it with the `REPLICATION` privilege when [Configuring the Database Cluster](#):

```
CREATE ROLE backup WITH LOGIN REPLICATION;
```

- If the `backup` role already exists, grant it with the `REPLICATION` privilege:

```
ALTER ROLE backup WITH REPLICATION;
```

- In the `pg_hba.conf` file, allow replication on behalf of the `backup` role.
- Make sure the parameter `max_wal_senders` is set high enough to leave at least one session available for the backup process.
- Set the parameter `wal_level` to be higher than `minimal`.

If you are planning to take PAGE backups in the STREAM mode or perform PITR with STREAM backups, you still have to configure WAL archiving, as explained in the section [Setting up continuous WAL archiving](#).

Once these steps are complete, you can start taking FULL, PAGE, DELTA, and PTRACK backups in the [STREAM](#) WAL mode.

Note

If you are planning to rely on `.pgpass` for authentication when running backup in STREAM mode, then `.pgpass` must contain credentials for `replication` database, used to establish connection via replication protocol. Example: `pghost:5432:replication:backup_user:my_strong_password`

Setting up Continuous WAL Archiving

Making backups in the PAGE backup mode, performing PITR and making backups with the [ARCHIVE](#) WAL delivery mode require [continuous WAL archiving](#) to be enabled. To set up continuous archiving in the cluster, complete the following steps:

- Make sure the `wal_level` parameter is higher than `minimal`.
- If you are configuring archiving on the primary, `archive_mode` must be set to `on` or `always`. To perform archiving on a standby, set this parameter to `always`.
- Set the `archive_command` parameter, as follows:

```
archive_command = '"install_dir/pg_probackup" archive-push -B "backup_dir" --  
instance=instance_name --wal-file-name=%f [remote_options]'
```

where `install_dir` is the installation directory of the `pg_probackup` version you are going to use, `backup_dir` and `instance_name` refer to the already initialized backup catalog instance for this database cluster, and `remote_options` only need to be specified to archive WAL on a remote host. For details about all possible `archive-push` parameters, see the section [archive-push](#).

Once these steps are complete, you can start making backups in the [ARCHIVE](#) WAL mode, backups in the PAGE backup mode, as well as perform PITR.

You can view the current state of the WAL archive using the `show` command. For details, see [the section called "Viewing WAL Archive Information"](#).

If you are planning to make PAGE backups and/or backups with [ARCHIVE](#) WAL mode from a standby server that generates a small amount of WAL traffic, without long waiting for WAL segment to fill up, consider setting the [archive_timeout](#) Postgres Pro parameter **on the primary**. The value of this parameter should be slightly lower than the `--archive-timeout` setting (5 minutes by default), so that there is enough time for the rotated segment to be streamed to a standby and sent to WAL archive before the backup is aborted because of `--archive-timeout`.

Note

Instead of using the [archive-push](#) command provided by `pg_probackup`, you can use any other tool to set up continuous archiving as long as it delivers WAL segments into `backup_dir/wal/instance_name` directory. If compression is used, it should be `gzip`, and `.gz` suffix in filename is mandatory.

Note

Instead of configuring continuous archiving by setting the `archive_mode` and `archive_command` parameters, you can opt for using the [pg_receivewal](#) utility. In this case, `pg_receivewal -D directory` option should point to `backup_dir/wal/instance_name` directory. `pg_probackup` supports WAL compression that can be done by `pg_receivewal`. “Zero Data Loss” archive strategy can be achieved only by using `pg_receivewal`.

Setting up Backups from a Standby

`pg_probackup` can take backups from a standby server. This requires the following additional setup:

- On the standby server, set the [hot_standby](#) parameter to `on`.
- On the primary server, set the [full_page_writes](#) parameter to `on`.
- To perform standalone backups on a standby, complete all the steps in section [Setting up STREAM Backups](#).
- To perform archive backups on a standby, complete all steps in section [Setting up continuous WAL archiving](#).

Once these steps are complete, you can start taking FULL, PAGE, DELTA, or PTRACK backups with appropriate WAL delivery mode: ARCHIVE or STREAM, from the standby server.

Backups from a standby server have the following limitations:

- If a standby is promoted to a primary during the backup process, the backup fails.
- All WAL records required for the backup must contain sufficient full-page writes. This requires you to enable `full_page_writes` on the primary, and not to use tools like `pg_compresslog` as [archive_command](#) to remove full-page writes from WAL files.

Setting up Cluster Verification

Logical verification of a database cluster requires the following additional setup. Role `backup` is used as an example:

- Install the [amcheck](#) or [amcheck_next](#) extension **in every database** of the cluster:

```
CREATE EXTENSION amcheck;
```

- Grant the following permissions to the `backup` role **in every database** of the cluster:

```
GRANT SELECT ON TABLE pg_catalog.pg_am TO backup;
GRANT SELECT ON TABLE pg_catalog.pg_class TO backup;
GRANT SELECT ON TABLE pg_catalog.pg_database TO backup;
GRANT SELECT ON TABLE pg_catalog.pg_namespace TO backup;
```

```
GRANT SELECT ON TABLE pg_catalog.pg_extension TO backup;
GRANT EXECUTE ON FUNCTION bt_index_check(regclass) TO backup;
GRANT EXECUTE ON FUNCTION bt_index_check(regclass, bool) TO backup;
GRANT EXECUTE ON FUNCTION bt_index_check(regclass, bool, bool) TO backup;
```

Setting up Partial Restore

If you are planning to use partial restore, complete the following additional step:

- Grant the read-only access to `pg_catalog.pg_database` to the backup role only in the database **used for connection** to Postgres Pro server:

```
GRANT SELECT ON TABLE pg_catalog.pg_database TO backup;
```

Configuring the Remote Mode

`pg_probackup` supports the remote mode that allows you to perform backup, restore and WAL archiving operations remotely. In this mode, the backup catalog is stored on a local system, while Postgres Pro instance to backup and/or to restore is located on a remote system. Currently the only supported remote protocol is SSH.

Set up SSH

If you are going to use `pg_probackup` in remote mode via SSH, complete the following steps:

1. Install `pg_probackup` on both systems: `backup_host` and `postgres_host`.
2. For communication between the hosts set up a passwordless SSH connection between the `backup_user` user on `backup_host` and the `postgres` user on `postgres_host`:

```
backup_user@backup_host:~$ ssh-copy-id postgres@postgres_host
```

Where:

- `backup_host` is the system with *backup catalog*.
 - `postgres_host` is the system with the Postgres Pro cluster.
 - `backup_user` is the OS user on `backup_host` used to run `pg_probackup`.
 - `postgres` is the user on `postgres_host` under which Postgres Pro cluster processes are running. For Postgres Pro 11 or higher a more secure approach can be used thanks to [group access](#) feature.
3. If you are going to rely on [continuous WAL archiving](#), set up a passwordless SSH connection between the `postgres` user on `postgres_host` and the `backup` user on `backup_host`:
- ```
postgres@postgres_host:~$ ssh-copy-id backup_user@backup_host
```
4. Make sure `pg_probackup` on `postgres_host` can be located when a connection via SSH is made. For example, for Bash, you can modify `PATH` in `~/.bashrc` of the `postgres` user (above the line in `bashrc` that exits the script for non-interactive shells). Alternatively, for `pg_probackup` commands, specify the path to the directory containing the `pg_probackup` binary on `postgres_host` via the [--remote-path](#) option.

`pg_probackup` in the remote mode via SSH works as follows:

- Only the following commands can be launched in the remote mode: [add-instance](#), [backup](#), [restore](#), [delete](#), [catchup](#), [archive-push](#), and [archive-get](#).
- Operating in remote mode requires `pg_probackup` binary to be installed on both local and remote systems. The versions of local and remote binary must be the same.
- When started in the remote mode, the main `pg_probackup` process on the local system connects to the remote system via SSH and launches one or more agent processes on the remote system, which are called *remote agents*. The number of remote agents is equal to the `-j/--threads` setting.
- The main `pg_probackup` process uses remote agents to access remote files and transfer data between local and remote systems.
- Remote agents try to minimize the network traffic and the number of round-trips between hosts.



- The main process is usually started on *backup\_host* and connects to *postgres\_host*, but in case of `archive-push` and `archive-get` commands the main process is started on *postgres\_host* and connects to *backup\_host*.
- Once data transfer is complete, remote agents are terminated and SSH connections are closed.
- If an error condition is encountered by a remote agent, then all agents are terminated and error details are reported by the main `pg_probackup` process, which exits with an error.
- Compression is always done on *postgres\_host*, while decompression is always done on *backup\_host*.

### Note

You can impose [additional restrictions](#) on SSH settings to protect the system in the event of account compromise.

### Note

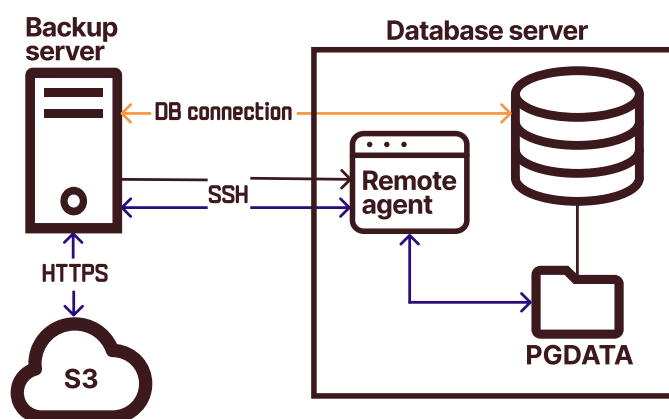
Setting the number of threads (`-j/--threads` option) to a value greater than 10 for `pg_probackup` working in the remote mode via SSH may result in the actual number of SSH connections exceeding the maximum allowed number of simultaneous SSH connections on the remote server and consequently lead to an “ERROR: Agent error: kex\_exchange\_identification: Connection closed by remote host” error. To correct the error, either reduce the number of `pg_probackup` threads or adjust the value of `MaxStartups` configuration parameter of the remote SSH server. If SSH is set up as a `xinetd` service on the remote server, adjust the value of the `xinetd per_source` configuration parameter rather than `MaxStartups`.

## Configuring S3 Connectivity

`pg_probackup` supports S3 interface for storing backups. Backup data is transferred to and from S3 without saving it in intermediate locations thus eliminating the need of having a large temporary storage.

An example configuration with a remote agent and a cloud storage (S3) is shown in [Figure H.1](#).

**Figure H.1. `pg_probackup` setup with a remote agent and S3**



In this figure, the following logical components are shown:

**Backup server**

A server where the main process of `pg_probackup` runs and where local backups are stored.

**Database server**

A server with a database instance that needs to be backed up or restored.

**Remote agent**

A secondary `pg_probackup` process running on the database server. Only applicable to the remote mode.

**Cloud storage**

A cloud storage for backups.

**Set up Access to S3 Storage**

If you are going to use `pg_probackup` with S3 interface, complete the following steps:

- Create a bucket with a unique and meaningful name in the S3 storage for you future backups.
- Create `ACCESS_KEY` and `SECRET_ACCESS_KEY` tokens to be used for secure connectivity instead of your username and password.
- For communication between `pg_probackup` and S3 server, set values of environment variables corresponding to your S3 server. For example:

```
export PG_PROBACKUP_S3_HOST=127.0.0.1
export PG_PROBACKUP_S3_PORT=9000
export PG_PROBACKUP_S3_REGION=ru-msk
export PG_PROBACKUP_S3_BUCKET_NAME=test1
export PG_PROBACKUP_S3_ACCESS_KEY=admin
export PG_PROBACKUP_S3_SECRET_ACCESS_KEY=password
export PG_PROBACKUP_S3_HTTPS=ON
```

Alternatively, you can provide S3 server settings in the S3 configuration file (see the `--s3-config-file` option in the section [S3 Options](#) for details).

It makes sense to specify S3 server settings if `--s3=minio`, as described in the section [S3 Options](#).

The following environment variables can be specified:

`PG_PROBACKUP_S3_HOST`

Address or list of addresses of the S3 server. A list of one or several semicolon-delimited addresses. Do not add a semicolon after the last address in the list. Each address can include the port number, separated by a colon. If the port number is not specified, the value of `PG_PROBACKUP_S3_PORT` is assumed. Do not add a colon if the port number is not specified.

For example:

```
export PG_PROBACKUP_S3_PORT=80
export PG_PROBACKUP_S3_HOST="127.0.0.1:9000;10.4.13.56:443;172.17.0.1"
```

In this example, for the “127.0.0.1” address, the port 9000 is explicitly specified, for “10.4.13.56”, the port 443 is specified, while for the “172.17.0.1” address, port 80, specified through `PG_PROBACKUP_S3_PORT`, will be used.

If any of the specified addresses gets unavailable while `pg_probackup` is in operation, requests to the S3 storage are distributed between the rest of the specified addresses. That is, when several addresses are specified, `pg_probackup` performs load balancing of S3 requests.

`PG_PROBACKUP_S3_PORT`

The port of the S3 server.

`PG_PROBACKUP_S3_REGION`

The region of the S3 server.

`PG_PROBACKUP_S3_BUCKET_NAME`

The name of the bucket on the S3 server.

PG\_PROBACKUP\_S3\_ACCESS\_KEY  
PG\_PROBACKUP\_S3\_SECRET\_ACCESS\_KEY

Secure tokens on the S3 server.

PG\_PROBACKUP\_S3\_HTTPS

The protocol to be used. Possible values:

- ON or HTTPS — use HTTPS
- Other than ON or HTTPS — use HTTP

PG\_PROBACKUP\_S3\_BUFFER\_SIZE

The size of the read/write buffer for communicating with S3, in MiB. The default is 16.

PG\_PROBACKUP\_S3\_RETRIES

The maximum number of attempts to execute an S3 request in case of failures. The default is 3.

PG\_PROBACKUP\_S3\_TIMEOUT

The maximum amount of time to execute an HTTP request to the S3 server, in seconds. The default is 300.

PG\_PROBACKUP\_S3\_IGNORE\_CERT\_VER

Don't verify the certificate host and peer. The default is ON.

PG\_PROBACKUP\_S3\_CA\_CERTIFICATE

Specify the path to file with trust Certificate Authority (CA) bundle.

PG\_PROBACKUP\_S3\_CA\_PATH

Specify the directory with trust CA certificates.

PG\_PROBACKUP\_S3\_CLIENT\_CERT

Setup SSL client certificate.

PG\_PROBACKUP\_S3\_CLIENT\_KEY

Setup private key file for TLS and SSL client certificate.

## Setting up PTRACK Backups

### Note

PTRACK versions lower than 2.0 are deprecated and not supported. Postgres Pro Standard and Postgres Pro Enterprise versions starting with 11.9.1 contain PTRACK 2.0. Upgrade your server to avoid issues in backups that you will take in future and be sure to take fresh backups of your clusters with the upgraded PTRACK since the backups taken with PTRACK 1.x might be corrupt.

If you are going to use PTRACK backups, complete the following additional steps.

### Note

The permissions required for the role that will perform PTRACK backups (the `backup` role in the examples below) are listed in [the section called “Configuring the Database Cluster”](#). The role must have permissions only in the database used for connection to the Postgres Pro server.

For Postgres Pro 11 or higher:

1. Create PTRACK extension:

```
CREATE EXTENSION ptrack;
```

2. To enable tracking page updates, set `ptrack.map_size` parameter to a positive integer and restart the server.

For optimal performance, it is recommended to set `ptrack.map_size` to  $N / 1024$ , where  $N$  is the size of the Postgres Pro cluster, in MB. If you set this parameter to a lower value, PTRACK is more likely to map several blocks together, which leads to false-positive results when tracking changed blocks and increases the incremental backup size as unchanged blocks can also be copied into the incremental backup. Setting `ptrack.map_size` to a higher value does not affect PTRACK operation, but it is not recommended to set this parameter to a value higher than 1024.

### Note

If you change the `ptrack.map_size` parameter value, the previously created PTRACK map file is cleared, and tracking newly changed blocks starts from scratch. Thus, you have to retake a full backup before taking incremental PTRACK backups after changing `ptrack.map_size`.

## Usage

### Creating a Backup

To create a backup, run the following command:

```
pg_probackup backup -B backup_dir --instance=instance_name -b backup_mode
```

Where *backup\_mode* can take one of the following values: [FULL](#), [DELTA](#), [PAGE](#), and [PTRACK](#).

When restoring a cluster from an incremental backup, `pg_probackup` relies on the parent full backup and all the incremental backups between them, which is called “the backup chain”. Thus, to perform incremental backups, it is necessary to have the last parent full backup with the status `OK` or `DONE` in the directory. If the parent full backup has the `MERGING` or `MERGED` status, an incremental backup cannot be performed.

For example, if merge has already been launched with a single full backup, an attempt to perform an incremental backup will end with the following messages:

```
WARNING: Valid full backup on current timeline 1 is not found, trying to look up on
previous timelines
WARNING: Cannot find valid backup on previous timelines
ERROR: Create new full backup before an incremental one
```

### ARCHIVE Mode

ARCHIVE is the default WAL delivery mode.

For example, to make a FULL backup in ARCHIVE mode, run:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL
```

ARCHIVE backups rely on [continuous archiving](#) to get WAL segments required to restore the cluster to a consistent state at the time the backup was taken.

When a backup is taken, `pg_probackup` ensures that WAL files containing WAL records between `Start LSN` and `Stop LSN` actually exist in `backup_dir/wal/instance_name` directory. `pg_probackup` also ensures that WAL records between `Start LSN` and `Stop LSN` can be parsed. This precaution eliminates the risk of silent WAL corruption.

## STREAM Mode

STREAM is the optional WAL delivery mode.

For example, to make a FULL backup in the STREAM mode, add the `--stream` flag to the command from the previous example:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --stream --temp-slot
```

The optional `--temp-slot` flag ensures that the required segments remain available if the WAL is rotated before the backup is complete.

Unlike backups in ARCHIVE mode, STREAM backups include all the WAL segments required to restore the cluster to a consistent state at the time the backup was taken.

During [backup](#) `pg_probackup` streams WAL files containing WAL records between `Start LSN` and `Stop LSN` to `backup_dir/backups/instance_name/backup_id/database/pg_wal` directory. To eliminate the risk of silent WAL corruption, `pg_probackup` also checks that WAL records between `Start LSN` and `Stop LSN` can be parsed.

Even if you are using [continuous archiving](#), STREAM backups can still be useful in the following cases:

- STREAM backups can be restored on the server that has no file access to WAL archive.
- STREAM backups enable you to restore the cluster state at the point in time for which WAL files in archive are no longer available.
- Backup in STREAM mode can be taken from a standby of a server that generates small amount of WAL traffic, without long waiting for WAL segment to fill up.

## Page Validation

If [data checksums](#) are enabled in the database cluster, `pg_probackup` uses this information to check correctness of data files during backup. While reading each page, `pg_probackup` checks whether the calculated checksum coincides with the checksum stored in the page header. This guarantees that the Postgres Pro instance and the backup itself have no corrupt pages. Note that `pg_probackup` reads database files directly from the filesystem, so under heavy write load during backup it can show false-positive checksum mismatches because of partial writes. If a page checksum mismatch occurs, the page is re-read and checksum comparison is repeated.

A page is considered corrupt if checksum comparison has failed more than 300 times. In this case, the backup is aborted.

Even if data checksums are not enabled, `pg_probackup` always performs sanity checks for page headers.

## External Directories

To back up a directory located outside of the data directory, use the optional `--external-dirs` parameter that specifies the path to this directory. If you would like to add more than one external directory, you can provide several paths separated by colons on Linux systems or semicolons on Windows systems.

For example, to include `/etc/dir1` and `/etc/dir2` directories into the full backup of your `instance_name` instance that will be stored under the `backup_dir` directory on Linux, run:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --external-dirs=/etc/dir1:/etc/dir2
```

Similarly, to include `C:\dir1` and `C:\dir2` directories into the full backup on Windows, run:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --external-dirs=C:\dir1;C:\dir2
```

`pg_probackup` recursively copies the contents of each external directory into a separate subdirectory in the backup catalog. Since external directories included into different backups do not have to be the same, when you are restoring the cluster from an incremental backup, only those directories that belong to this particular backup will be restored. Any external directories stored in the previous backups will be ignored.

To include the same directories into each backup of your instance, you can specify them in the `pg_probackup.conf` configuration file using the [set-config](#) command with the `--external-dirs` option.

## Performing Cluster Verification

To verify that Postgres Pro database cluster is not corrupt, run the following command:

```
pg_probackup checkdb [-B backup_dir [--instance=instance_name]] [-D data_dir]
[connection_options]
```

This command performs physical verification of all data files located in the specified data directory by running page header sanity checks, as well as block-level checksum verification if checksums are enabled. If a corrupt page is detected, `checkdb` continues cluster verification until all pages in the cluster are validated.

By default, similar [page validation](#) is performed automatically while a backup is taken by `pg_probackup`. The `checkdb` command enables you to perform such page validation on demand, without taking any backup copies, even if the cluster is not backed up using `pg_probackup` at all.

To perform cluster verification, `pg_probackup` needs to connect to the cluster to be verified. In general, it is enough to specify the backup instance of this cluster for `pg_probackup` to determine the required connection options. However, if `-B` and `--instance` options are omitted, you have to provide [connection options](#) and `data_dir` via environment variables or command-line options.

Physical verification cannot detect logical inconsistencies, missing or nullified blocks and entire files, or similar anomalies. Extensions [amcheck](#) and [amcheck\\_next](#) provide a partial solution to these problems.

If you would like, in addition to physical verification, to verify all indexes in all databases using these extensions, you can specify the `--amcheck` flag when running the [checkdb](#) command:

```
pg_probackup checkdb -D data_dir --amcheck [connection_options]
```

You can skip physical verification by specifying the `--skip-block-validation` flag. In this case, you can omit `backup_dir` and `data_dir` options, only [connection options](#) are mandatory:

```
pg_probackup checkdb --amcheck --skip-block-validation [connection_options]
```

Logical verification can be done more thoroughly with the `--heapallindexed` flag by checking that all heap tuples that should be indexed are actually indexed, but at the higher cost of CPU, memory, and I/O consumption.

## Validating a Backup

`pg_probackup` calculates checksums for each file in a backup during the backup process. The process of checking checksums of backup data files is called *the backup validation*. By default, validation is run immediately after the backup is taken and right before the restore, to detect possible backup corruption.

### Note

The backup validation includes checking checksums for CFS files.

If you would like to skip backup validation, you can specify the `--no-validate` flag when running [backup](#) and [restore](#) commands.

To ensure that all the required backup files are present and can be used to restore the database cluster, you can run the [validate](#) command with the exact [recovery target options](#) you are going to use for recovery.

For example, to check that you can restore the database cluster from a backup copy up to transaction ID 4242, run this command:

```
pg_probackup validate -B backup_dir --instance=instance_name --recovery-target-xid=4242
```

If validation completes successfully, `pg_probackup` displays the corresponding message. If validation fails, you will receive an error message with the exact time, transaction ID, and LSN up to which the recovery is possible.

If you specify `backup_id` via `-i/--backup-id` option, then only the backup copy with specified backup ID will be validated. If `backup_id` is specified with [recovery target options](#), the `validate` command will check whether it is possible to restore the specified backup to the specified recovery target.

For example, to check that you can restore the database cluster from a backup copy with the SBOL6P backup ID up to the specified timestamp, run this command:

```
pg_probackup validate -B backup_dir --instance=instance_name -i SBOL6P --recovery-
target-time="2024-04-10 18:18:26+03"
```

If you specify the `backup_id` of an incremental backup, all its parents starting from FULL backup will be validated.

If you omit all the parameters, all backups are validated.

## Restoring a Cluster

To restore the database cluster from a backup, run the [restore](#) command with at least the following options:

```
pg_probackup restore -B backup_dir --instance=instance_name -i backup_id
```

Where:

- `backup_dir` is the backup catalog that stores all backup files and meta information.
- `instance_name` is the backup instance for the cluster to be restored.
- `backup_id` specifies the backup to restore the cluster from. If you omit this option, `pg_probackup` uses the latest valid backup available for the specified instance. If you specify an incremental backup to restore, `pg_probackup` automatically restores the underlying full backup and then sequentially applies all the necessary increments.

Once the `restore` command is complete, start the database service.

If you restore [ARCHIVE](#) backups, perform [PITR](#), or specify the `--restore-as-replica` flag with the `restore` command to set up a standby server, `pg_probackup` creates a recovery configuration file once all data files are copied into the target directory. This file includes the minimal settings required for recovery, except for the password in the [primary\\_conninfo](#) parameter; you have to add the password manually or use the `--primary-conninfo` option, if required. For Postgres Pro 11, recovery settings are written into the `recovery.conf` file. Starting from Postgres Pro 12, `pg_probackup` writes these settings into the `probackup_recovery.conf` file and then includes it into `postgresql.auto.conf`.

If you are restoring a [STREAM](#) backup, the restore is complete at once, with the cluster returned to a self-consistent state at the point when the backup was taken. For [ARCHIVE](#) backups, Postgres Pro replays all available archived WAL segments, so the cluster is restored to the latest state possible within the current timeline. You can change this behavior by using the [recovery target options](#) with the `restore` command, as explained in [the section called “Performing Point-in-Time \(PITR\) Recovery”](#).

If the cluster to restore contains tablespaces, `pg_probackup` restores them to their original location by default. To restore tablespaces to a different location, use the `--tablespace-mapping/-T` option. Otherwise, restoring the cluster on the same host will fail if tablespaces are in use, because the backup would have to be written to the same directories.

When using the `--tablespace-mapping/-T` option, you must provide absolute paths to the old and new tablespace directories. If a path happens to contain an equals sign (=), escape it with a backslash. This option can be specified multiple times for multiple tablespaces. For example:

```
pg_probackup restore -B backup_dir --instance=instance_name -D data_dir -j 4 -
i backup_id -T tablespace1_dir=tablespace1_newdir -T tablespace2_dir=tablespace2_newdir
```



To restore the cluster on a remote host, follow the instructions in [the section called “Using pg\\_probackup in the Remote Mode”](#).

**Note**

By default, the [restore](#) command validates the specified backup before restoring the cluster. If you run regular backup validations and would like to save time when restoring the cluster, you can specify the `--no-validate` flag to skip validation and speed up the recovery.

**Incremental Restore**

The speed of restore from backup can be significantly improved by replacing only invalid and changed pages in already existing Postgres Pro data directory using [incremental restore options](#) with the [restore](#) command.

To restore the database cluster from a backup in incremental mode, run the [restore](#) command with the following options:

```
pg_probackup restore -B backup_dir --instance=instance_name -D data_dir -
I incremental_mode
```

Where *incremental\_mode* can take one of the following values:

- CHECKSUM — read all data files in the data directory, validate header and checksum in every page and replace only invalid pages and those with checksum and LSN not matching with corresponding page in backup. This is the simplest, the most fool-proof incremental mode. Recommended to use by default.
- LSN — read the *pg\_control* in the data directory to obtain redo LSN and redo TLI, which allows you to determine a point in history(shiftpoint), where data directory state shifted from target back-up chain history. If shiftpoint is not within reach of backup chain history, then restore is aborted. If shiftpoint is within reach of backup chain history, then read all data files in the data directory, validate header and checksum in every page and replace only invalid pages and those with LSN greater than shiftpoint. This mode offers a greater speed up compared to CHECKSUM, but rely on two conditions to be met. First, [data\\_checksums](#) parameter must be enabled in data directory (to avoid corruption due to hint bits). This condition will be checked at the start of incremental restore and the operation will be aborted if checksums are disabled. Second, the *pg\_control* file must be synched with state of data directory. This condition cannot be checked at the start of restore, so it is a user responsibility to ensure that *pg\_control* contain valid information. Therefore it is not recommended to use LSN mode in any situation, where *pg\_control* cannot be trusted or has been tampered with: after *pg\_resetxlog* execution, after restore from backup without recovery been run, etc.
- NONE — regular restore without any incremental optimizations.

Regardless of chosen incremental mode, *pg\_probackup* will check, that postmaster in given destination directory is not running and *system-identifier* is the same as in the backup.

Suppose you want to return an old primary as a replica after switchover using incremental restore in LSN mode:

| =====    |         |      |               |                               |            |            |          |     |      |
|----------|---------|------|---------------|-------------------------------|------------|------------|----------|-----|------|
| Instance | Version | ID   | Recovery Time |                               | Mode       |            | WAL Mode | TLI | Time |
| Data     | WAL     | Zalg | Zratio        | Start LSN                     | Stop LSN   | Status     |          |     |      |
| =====    |         |      |               |                               |            |            |          |     |      |
| node     | 17      |      | SBOL8S        | 2024-04-09 18:19:43.707720+03 | DELTA      | STREAM     | 16/15    |     |      |
| 3s       | 114MB   | 64MB | lz4           | 1.42                          | 0/3C003020 | 0/3E8D4930 | OK       |     |      |
| node     | 17      |      | SBOL8G        | 2024-04-09 18:19:32.594670+03 | PTRACK     | STREAM     | 15/15    |     |      |
| 4s       | 30MB    | 16MB | zlib          | 2.23                          | 0/31000028 | 0/310029E0 | OK       |     |      |
| node     | 17      |      | SBOL83        | 2024-04-09 18:19:22.269595+03 | PAGE       | STREAM     | 15/15    |     |      |
| 7s       | 46MB    | 32MB | pglz          | 1.44                          | 0/29000028 | 0/2A0000F8 | OK       |     |      |



```
node 17 SBOL7P 2024-04-09 18:19:06.557301+03 FULL STREAM 15/0
6s 144MB 16MB zstd 2.47 0/22000028 0/220001C8 OK
```

```
backup_user@backup_host:~$ pg_probackup restore -B /mnt/backups --instance=node -R -I
lsn
INFO: Destination directory and tablespace directories are empty, disable incremental
restore
INFO: Validating parents for backup SBOL8S
INFO: Validating backup SBOL7P
INFO: Backup SBOL7P data files are valid
INFO: Validating backup SBOL83
INFO: Backup SBOL83 data files are valid
INFO: Validating backup SBOL8G
INFO: Backup SBOL8G data files are valid
INFO: Validating backup SBOL8S
INFO: Backup SBOL8S data files are valid
INFO: Backup SBOL8S WAL segments are valid
INFO: Backup SBOL8S is valid.
INFO: Restoring the database from the backup starting at 2024-04-09 18:19:40+03
INFO: Start restoring backup files. PGDATA size: 616MB
INFO: Backup files are restored. Transferred bytes: 616MB, time elapsed: 2s
INFO: Restore incremental ratio (less is better): 100% (616MB/616MB)
INFO: Syncing restored files to disk
INFO: Restored backup files are synced, time elapsed: 2s
INFO: Restore of backup SBOL8S completed.
```

### Note

Incremental restore is possible only for backups with `program_version` equal or greater than 2.4.0.

### Partial Restore

If you have enabled [partial restore](#) before taking backups, you can restore only some of the databases using [partial restore options](#) with the [restore](#) commands.

To restore the specified databases only, run the [restore](#) command with the following options:

```
pg_probackup restore -B backup_dir --instance=instance_name --db-include=database_name
```

The `--db-include` option can be specified multiple times. For example, to restore only databases `db1` and `db2`, run the following command:

```
pg_probackup restore -B backup_dir --instance=instance_name --db-include=db1 --db-include=db2
```

To exclude one or more databases from restore, use the `--db-exclude` option:

```
pg_probackup restore -B backup_dir --instance=instance_name --db-exclude=database_name
```

The `--db-exclude` option can be specified multiple times. For example, to exclude the databases `db1` and `db2` from restore, run the following command:

```
pg_probackup restore -B backup_dir --instance=instance_name --db-exclude=db1 --db-exclude=db2
```

Partial restore relies on lax behavior of Postgres Pro recovery process toward truncated files. For recovery to work properly, files of excluded databases are restored as files of zero size. After the Postgres Pro cluster is successfully started, you must drop the excluded databases using `DROP DATABASE` command.

To decouple a single cluster containing multiple databases into separate clusters with minimal downtime, you can do partial restore of the cluster as a standby using the `--restore-as-replica` option for specific databases.

### Note

The `template0` and `template1` databases are always restored.

### Note

Due to recovery specifics of Postgres Pro versions earlier than 12, it is advisable that you set the `hot_standby` parameter to `off` when running partial restore of a Postgres Pro cluster of version earlier than 12. Otherwise the recovery may fail.

## Performing Point-in-Time (PITR) Recovery

If you have enabled [continuous WAL archiving](#) before taking backups, you can restore the cluster to its state at an arbitrary point in time (recovery target) using [recovery target options](#) with the `restore` command.

You can use both STREAM and ARCHIVE backups for point in time recovery as long as the WAL archive is available at least starting from the time the backup was taken. If `-i/--backup-id` option is omitted, `pg_probackup` automatically chooses the backup that is the closest to the specified recovery target and starts the restore process, otherwise `pg_probackup` will try to restore the specified backup to the specified recovery target.

- To restore the cluster state at the exact time, specify the `--recovery-target-time` option, in the timestamp format. For example:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target-time="2024-04-10 18:18:26+03"
```

- To restore the cluster state up to a specific transaction ID, use the `--recovery-target-xid` option:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target-xid=687
```

- To restore the cluster state up to the specific LSN, use `--recovery-target-lsn` option:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target-lsn=16/B374D848
```

- To restore the cluster state up to the specific named restore point, use `--recovery-target-name` option:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target-name="before_app_upgrade"
```

- To restore the backup to the latest state available in the WAL archive, use `--recovery-target` option with `latest` value:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target="latest"
```

- To restore the cluster to the earliest point of consistency, use `--recovery-target` option with the `immediate` value:

```
pg_probackup restore -B backup_dir --instance=instance_name --recovery-target='immediate'
```

## Using pg\_probackup in the Remote Mode

pg\_probackup supports the remote mode that allows you to perform `backup` and `restore` operations remotely via SSH. In this mode, the backup catalog is stored on a local system, while Postgres Pro instance to be backed up is located on a remote system. You must have pg\_probackup installed on both systems.

### Note

pg\_probackup relies on passwordless SSH connection for communication between the hosts.

### Note

In addition to SSH connection, pg\_probackup uses a regular connection to the database to manage the remote operation. See the section [Configuring the Database Cluster](#) for details of how to set up a database connection.

The typical workflow is as follows:

- On your backup host, configure pg\_probackup as explained in the section [Installation and Setup](#). For the `add-instance` and `set-config` commands, make sure to specify `remote options` that point to the database host with the Postgres Pro instance.
- If you would like to take remote backups in `PAGE` mode, or rely on `ARCHIVE` WAL delivery mode, or use `PITR`, configure continuous WAL archiving from the database host to the backup host as explained in the section [Setting up continuous WAL archiving](#). For the `archive-push` and `archive-get` commands, you must specify the `remote options` that point to the backup host with the backup catalog.
- Run `backup` or `restore` commands with `remote options on the backup host`. pg\_probackup connects to the remote system via SSH and creates a backup locally or restores the previously taken backup on the remote system, respectively.

For example, to create an archive full backup of a Postgres Pro cluster located on a remote system with host address `192.168.0.2` on behalf of the `postgres` user via SSH connection through port 2302, run:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --remote-user=postgres --remote-host=192.168.0.2 --remote-port=2302
```

To restore the latest available backup on a remote system with host address `192.168.0.2` on behalf of the `postgres` user via SSH connection through port 2302, run:

```
pg_probackup restore -B backup_dir --instance=instance_name --remote-user=postgres --remote-host=192.168.0.2 --remote-port=2302
```

Restoring an `ARCHIVE` backup or performing `PITR` in the remote mode require additional information: destination address, port and username for establishing an SSH connection **from** the host with database **to** the host with the backup catalog. This information will be used by the `restore_command` to copy WAL segments from the archive to the Postgres Pro `pg_wal` directory.

To solve this problem, you can use [Remote WAL Archive Options](#).

For example, to restore latest backup on remote system using remote mode through SSH connection to user `postgres` on host with address `192.168.0.2` via port 2302 and user `backup` on backup catalog host with address `192.168.0.3` via port 2303, run:

```
pg_probackup restore -B backup_dir --instance=instance_name --remote-user=postgres --remote-host=192.168.0.2 --remote-port=2302 --archive-host=192.168.0.3 --archive-port=2303 --archive-user=backup
```

Provided arguments will be used to construct the `restore_command`:

```
restore_command = '"install_dir/pg_probackup" archive-get -B "backup_dir" --
instance=instance_name --wal-file-path=%p --wal-file-name=%f --remote-host=192.168.0.3
--remote-port=2303 --remote-user=backup'
```

Alternatively, you can use the `--restore-command` option to provide the entire `restore_command`:

```
pg_probackup restore -B backup_dir --instance=instance_name --remote-user=postgres
--remote-host=192.168.0.2 --remote-port=2302 --restore-command='"install_dir/
pg_probackup" archive-get -B "backup_dir" --instance=instance_name --wal-file-path=%p
--wal-file-name=%f --remote-host=192.168.0.3 --remote-port=2303 --remote-user=backup'
```

### Note

The remote mode is currently unavailable for Windows systems.

## Running pg\_probackup on Parallel Threads

[backup](#), [restore](#), [merge](#), [delete](#), [catchup](#), [checkdb](#), and [validate](#) processes can be executed on several parallel threads. This can significantly speed up `pg_probackup` operation given enough resources (CPU cores, disk, and network bandwidth).

Parallel execution is controlled by the `-j/--threads` command-line option. For example, to create a backup using four parallel threads, run:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL -j 4
```

### Note

Parallel restore applies only to copying data from the backup catalog to the data directory of the cluster. When Postgres Pro server is started, WAL records need to be replayed, and this cannot be done in parallel.

## Configuring pg\_probackup

Once the backup catalog is initialized and a new backup instance is added, you can use the `pg_probackup.conf` configuration file located in the `backup_dir/backups/instance_name` directory to fine-tune `pg_probackup` configuration.

For example, [backup](#) and [checkdb](#) commands use a regular Postgres Pro connection. To avoid specifying [connection options](#) each time on the command line, you can set them in the `pg_probackup.conf` configuration file using the [set-config](#) command.

### Note

It is **not recommended** to edit `pg_probackup.conf` manually.

Initially, `pg_probackup.conf` contains the following settings:

- `PGDATA` — the path to the data directory of the cluster to back up.
- `system-identifier` — the unique identifier of the Postgres Pro instance.

Additionally, you can define [remote](#), [retention](#), [logging](#), and [compression](#) settings using the `set-config` command:

```
pg_probackup set-config -B backup_dir --instance=instance_name
```

```
[--external-dirs=external_directory_path] [remote_options] [connection_options]
[retention_options] [logging_options]
```

To view the current settings, run the following command:

```
pg_probackup show-config -B backup_dir --instance=instance_name
```

You can override the settings defined in `pg_probackup.conf` when running `pg_probackup` [commands](#) via the corresponding environment variables and/or command line options.

### Specifying Connection Settings

If you define connection settings in the `pg_probackup.conf` configuration file, you can omit connection options in all the subsequent `pg_probackup` commands. However, if the corresponding environment variables are set, they get higher priority. The options provided on the command line overwrite both environment variables and configuration file settings.

If nothing is given, the default values are taken. By default `pg_probackup` tries to use local connection via Unix domain socket (`localhost` on Windows) and tries to get the database name and the user name from the `PGUSER` environment variable or the current OS user name.

### Managing the Backup Catalog

With `pg_probackup`, you can manage backups from the command line:

- [View backup information](#)
- [View WAL Archive Information](#)
- [Validate backups](#)
- [Merge backups](#)
- [Delete backups](#)

### Viewing Backup Information

To view the list of existing backups for every instance, run the command:

```
pg_probackup show -B backup_dir
```

`pg_probackup` displays the list of all the available backups. For example:

```
BACKUP INSTANCE 'node'
```

| =====    |         |      |               |            |                    |        |          |     |      |
|----------|---------|------|---------------|------------|--------------------|--------|----------|-----|------|
| Instance | Version | ID   | Recovery Time |            | Mode               |        | WAL Mode | TLI | Time |
| Data     | WAL     | Zalg | Zratio        | Start LSN  | Stop LSN           | Status |          |     |      |
| =====    |         |      |               |            |                    |        |          |     |      |
| node     | 17      |      | SBOL94        | 2024-04-09 | 18:19:56.603355+03 | FULL   | ARCHIVE  | 1/0 | 6s   |
| 377MB    | 16MB    | lz4  | 1.46          | 0/41000028 | 0/420000C0         | OK     |          |     |      |
| node     | 17      |      | SBOL8S        | 2024-04-09 | 18:19:43.707720+03 | DELTA  | STREAM   | 1/1 | 3s   |
| 114MB    | 64MB    | lz4  | 1.42          | 0/3C003020 | 0/3E8D4930         | OK     |          |     |      |
| node     | 17      |      | SBOL8G        | 2024-04-09 | 18:19:32.594670+03 | PTRACK | STREAM   | 1/1 | 4s   |
| 30MB     | 16MB    | zlib | 2.23          | 0/31000028 | 0/310029E0         | OK     |          |     |      |
| node     | 17      |      | SBOL83        | 2024-04-09 | 18:19:22.269595+03 | PAGE   | STREAM   | 1/1 | 7s   |
| 46MB     | 32MB    | pglz | 1.44          | 0/29000028 | 0/2A0000F8         | OK     |          |     |      |
| node     | 17      |      | SBOL7P        | 2024-04-09 | 18:19:06.557301+03 | FULL   | STREAM   | 1/0 | 6s   |
| 144MB    | 16MB    | zstd | 2.47          | 0/22000028 | 0/220001C8         | OK     |          |     |      |

For each backup, the following information is provided:

- Instance — the instance name.
- Version — Postgres Pro major version.
- ID — the backup identifier.
- Recovery time — the earliest moment for which you can restore the state of the database cluster.
- Mode — the method used to take this backup. Possible values: FULL, PAGE, DELTA, PTRACK.

- **WAL Mode** — WAL delivery mode. Possible values: `STREAM` and `ARCHIVE`.
- **TLI** — timeline identifiers of the current backup and its parent.
- **Time** — the time it took to perform the backup.
- **Data** — the size of the data files in this backup. This value does not include the size of WAL files. For `STREAM` backups, the total size of the backup can be calculated as `Data + WAL`.
- **WAL** — the uncompressed size of WAL files that need to be applied during recovery for the backup to reach a consistent state.
- **compress-alg** — compression algorithm used during backup. Possible values: `zlib`, `pglz`, `lz4`, `zstd`, `none`.
- **Zratio** — compression ratio calculated as “uncompressed-bytes” / “data-bytes”.
- **Start LSN** — WAL log sequence number corresponding to the start of the backup process. REDO point for Postgres Pro recovery process to start from.
- **Stop LSN** — WAL log sequence number corresponding to the end of the backup process. Consistency point for Postgres Pro recovery process.
- **Status** — backup status. Possible values:
  - **OK** — the backup is complete and valid.
  - **DONE** — the backup is complete, but was not validated.
  - **RUNNING** — the backup is in progress.
  - **MERGING** — the backup is being merged.
  - **MERGED** — the backup data files were successfully merged, but its metadata is in the process of being updated. Only full backups can have this status.
  - **DELETING** — the backup files are being deleted.
  - **CORRUPT** — some of the backup files are corrupt.
  - **ERROR** — the backup was aborted because of an unexpected error.
  - **ORPHAN** — the backup is invalid because one of its parent backups is corrupt or missing.
  - **HIDDEN\_FOR\_TEST** — a test script marked the backup as nonexistent. (`pg_probackup` never sets this status by itself.)

You can restore the cluster from the backup only if the backup status is `OK` or `DONE`.

To get more detailed information about the backup, run the `show` command with the backup ID:

```
pg_probackup show -B backup_dir --instance=instance_name -i backup_id
```

The sample output is as follows:

```
#Configuration
backup-mode = FULL
stream = false
compress-alg = lz4
compress-level = 1
from-replica = false

#Compatibility
block-size = 8192
xlog-block-size = 8192
checksum-version = 1
program-version = 2.7.3
server-version = 17

#Result backup info
timelineid = 1
start-lsn = 0/41000028
stop-lsn = 0/420000C0
```

```
start-time = '2024-04-09 18:19:52+03'
end-time = '2024-04-09 18:19:58+03'
end-validation-time = '2024-04-09 18:19:59+03'
recovery-xid = 757
recovery-time = '2024-04-09 18:19:56.603355+03'
data-bytes = 395651278
wal-bytes = 16777216
uncompressed-bytes = 578552566
pgdata-bytes = 578552248
status = OK
primary_conninfo = 'user=backup channel_binding=prefer host=localhost
port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLSv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable'
content-crc = 3862224379
```

Detailed output has additional attributes:

- `compress-alg` — compression algorithm used during backup. Possible values: `zlib`, `pglz`, `lz4`, `zstd`, `none`.
- `compress-level` — compression level used during backup.
- `from-replica` — was this backup taken on a standby? Possible values: `1`, `0`.
- `block-size` — the [block\\_size](#) setting of Postgres Pro cluster at the backup start.
- `checksum-version` — are [data\\_checksums](#) enabled in the backed up Postgres Pro cluster? Possible values: `1`, `0`.
- `program-version` — full version of `pg_probackup` binary used to create the backup.
- `start-time` — the backup start time.
- `end-time` — the backup end time.
- `end-validation-time` — the backup validation end time.
- `expire-time` — the point in time when a pinned backup can be removed in accordance with retention policy. This attribute is only available for pinned backups.
- `uncompressed-bytes` — the size of data files before adding page headers and applying compression. You can evaluate the effectiveness of compression by comparing `uncompressed-bytes` to `data-bytes` if compression is used.
- `pgdata-bytes` — the size of Postgres Pro cluster data files at the time of backup. You can evaluate the effectiveness of an incremental backup by comparing `pgdata-bytes` to `uncompressed-bytes`.
- `recovery-xid` — transaction ID at the backup end time.
- `parent-backup-id` — ID of the parent backup. Available only for incremental backups.
- `primary_conninfo` — libpq connection parameters used to connect to the Postgres Pro cluster to take this backup. The password is not included.
- `note` — text note attached to backup.
- `content-crc` — CRC32 checksum of `backup_content.control` file. It is used to detect corruption of backup metainformation.

You can also get the detailed information about the backup in the JSON format:

```
pg_probackup show -B backup_dir --instance=instance_name --format=json -i backup_id
```

The sample output is as follows:

```
[
 {
 "instance": "node",
 "backups": [
 {
 "id": "SBOL94",
 "status": "OK",
 "start-time": "2024-04-09 18:19:52+03",
 "backup-mode": "FULL",
```

```

 "wal": "ARCHIVE",
 "compress-alg": "lz4",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,
 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 16,
 "parent-tli": 2,
 "start-lsn": "0/41000028",
 "stop-lsn": "0/420000C0",
 "end-time": "2024-04-09 18:19:58+03",
 "end-validation-time": "2024-04-09 18:19:59+03",
 "recovery-xid": 757,
 "recovery-time": "2024-04-09 18:19:56.603355+03",
 "data-bytes": 395651278,
 "wal-bytes": 16777216,
 "uncompressed-bytes": 578552566,
 "pgdata-bytes": 578552248,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLShv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 3862224379
 }
}
]

```

## Viewing WAL Archive Information

To view the information about WAL archive for every instance, run the command:

```
pg_probackup show -B backup_dir [--instance=instance_name] --archive
```

**pg\_probackup** displays the list of all the available WAL files grouped by timelines. For example:

```

INFO: checking WAL file name "000000010000000000000001B"
INFO: checking WAL file name "000000010000000000000001C"
INFO: checking WAL file name "000000010000000000000001D"
INFO: checking WAL file name "000000010000000000000001E"
INFO: checking WAL file name "000000010000000000000001F"
INFO: checking WAL file name "0000000100000000000000020"
INFO: checking WAL file name "0000000100000000000000021"
INFO: checking WAL file name "0000000100000000000000022"
INFO: checking WAL file name "0000000100000000000000022.00000028.backup"
INFO: checking WAL file name "0000000100000000000000023"
INFO: checking WAL file name "0000000100000000000000024"
INFO: checking WAL file name "0000000100000000000000025"
INFO: checking WAL file name "0000000100000000000000026"
INFO: checking WAL file name "0000000100000000000000027"
INFO: checking WAL file name "0000000100000000000000028"
INFO: checking WAL file name "0000000100000000000000029"
INFO: checking WAL file name "0000000100000000000000029.00000028.backup"
INFO: checking WAL file name "000000010000000000000002A"
INFO: checking WAL file name "000000010000000000000002B"
INFO: checking WAL file name "000000010000000000000002C"
INFO: checking WAL file name "000000010000000000000002D"

```



```

INFO: checking WAL file name "00000001000000000000000002E"
INFO: checking WAL file name "00000001000000000000000002F"
INFO: checking WAL file name "000000010000000000000000030"
INFO: checking WAL file name "000000010000000000000000031"
INFO: checking WAL file name "000000010000000000000000031.00000028.backup"
INFO: checking WAL file name "000000010000000000000000032"
INFO: checking WAL file name "000000010000000000000000033"
INFO: checking WAL file name "000000010000000000000000034"
INFO: checking WAL file name "000000010000000000000000035"
INFO: checking WAL file name "000000010000000000000000036"
INFO: checking WAL file name "000000010000000000000000037"
INFO: checking WAL file name "000000010000000000000000038"
INFO: checking WAL file name "000000010000000000000000039"
INFO: checking WAL file name "00000001000000000000000003A"
INFO: checking WAL file name "00000001000000000000000003B"
INFO: checking WAL file name "00000001000000000000000003C"
INFO: checking WAL file name "00000001000000000000000003C.00003020.backup"
INFO: checking WAL file name "00000001000000000000000003D"
INFO: checking WAL file name "00000001000000000000000003E"
INFO: checking WAL file name "00000001000000000000000003F"
INFO: checking WAL file name "000000010000000000000000040"
INFO: checking WAL file name "000000010000000000000000041"
INFO: checking WAL file name "000000010000000000000000041.00000028.backup"
INFO: checking WAL file name "000000010000000000000000042"

```

```

ARCHIVE INSTANCE 'node'

```

| TLI      | Parent TLI | Switchpoint | Min Segno                   | Max Segno                   | N  |
|----------|------------|-------------|-----------------------------|-----------------------------|----|
| segments | Size       | Zratio      | N backups                   | Status                      |    |
| 1        | 0          | 0/0         | 00000001000000000000000001B | 000000010000000000000000042 | 40 |
|          | 640MB      | 1.00        | 5                           | OK                          |    |

For each timeline, the following information is provided:

- **TLI** — timeline identifier.
- **Parent TLI** — identifier of the timeline from which this timeline branched off.
- **Switchpoint** — LSN of the moment when the timeline branched off from its parent timeline.
- **Min Segno** — the first WAL segment belonging to the timeline.
- **Max Segno** — the last WAL segment belonging to the timeline.
- **N segments** — number of WAL segments belonging to the timeline.
- **Size** — the size that files take on disk.
- **Zalg** — compression algorithm used during backup. Possible values: `zlib`, `pglz`, `lz4`, `zstd`, `none`.
- **Zratio** — compression ratio calculated as  $N \text{ segments} * wal\_segment\_size * wal\_block\_size / \text{Size}$ .
- **N backups** — number of backups belonging to the timeline. To get the details about backups, use the JSON format.
- **Status** — status of the WAL archive for this timeline. Possible values:
  - **OK** — all WAL segments between **Min Segno** and **Max Segno** are present.
  - **DEGRADED** — some WAL segments between **Min Segno** and **Max Segno** are missing. To find out which files are lost, view this report in the JSON format. This status may appear if several WAL files (in the middle of the sequence) were deleted by the `delete` command with the `--delete-wal` option according to the retention policy. This status does not affect the restore correctness, but it can be impossible to perform PITR of the cluster to some recovery targets.

To get more detailed information about the WAL archive in the JSON format, run the command:

```

pg_probackup show -B backup_dir [--instance=instance_name] --archive --format=json

```

The sample output is as follows:

```
INFO: checking WAL file name "0000000010000000000000000001B"
INFO: checking WAL file name "0000000010000000000000000001C"
INFO: checking WAL file name "0000000010000000000000000001D"
INFO: checking WAL file name "0000000010000000000000000001E"
INFO: checking WAL file name "0000000010000000000000000001F"
INFO: checking WAL file name "00000000100000000000000000020"
INFO: checking WAL file name "00000000100000000000000000021"
INFO: checking WAL file name "00000000100000000000000000022"
INFO: checking WAL file name "00000000100000000000000000022.00000028.backup"
INFO: checking WAL file name "00000000100000000000000000023"
INFO: checking WAL file name "00000000100000000000000000024"
INFO: checking WAL file name "00000000100000000000000000025"
INFO: checking WAL file name "00000000100000000000000000026"
INFO: checking WAL file name "00000000100000000000000000027"
INFO: checking WAL file name "00000000100000000000000000028"
INFO: checking WAL file name "00000000100000000000000000029"
INFO: checking WAL file name "00000000100000000000000000029.00000028.backup"
INFO: checking WAL file name "0000000010000000000000000002A"
INFO: checking WAL file name "0000000010000000000000000002B"
INFO: checking WAL file name "0000000010000000000000000002C"
INFO: checking WAL file name "0000000010000000000000000002D"
INFO: checking WAL file name "0000000010000000000000000002E"
INFO: checking WAL file name "0000000010000000000000000002F"
INFO: checking WAL file name "00000000100000000000000000030"
INFO: checking WAL file name "00000000100000000000000000031"
INFO: checking WAL file name "00000000100000000000000000031.00000028.backup"
INFO: checking WAL file name "00000000100000000000000000032"
INFO: checking WAL file name "00000000100000000000000000033"
INFO: checking WAL file name "00000000100000000000000000034"
INFO: checking WAL file name "00000000100000000000000000035"
INFO: checking WAL file name "00000000100000000000000000036"
INFO: checking WAL file name "00000000100000000000000000037"
INFO: checking WAL file name "00000000100000000000000000038"
INFO: checking WAL file name "00000000100000000000000000039"
INFO: checking WAL file name "0000000010000000000000000003A"
INFO: checking WAL file name "0000000010000000000000000003B"
INFO: checking WAL file name "0000000010000000000000000003C"
INFO: checking WAL file name "0000000010000000000000000003C.00003020.backup"
INFO: checking WAL file name "0000000010000000000000000003D"
INFO: checking WAL file name "0000000010000000000000000003E"
INFO: checking WAL file name "0000000010000000000000000003F"
INFO: checking WAL file name "00000000100000000000000000040"
INFO: checking WAL file name "00000000100000000000000000041"
INFO: checking WAL file name "00000000100000000000000000041.00000028.backup"
INFO: checking WAL file name "00000000100000000000000000042"
[
 {
 "instance": "node",
 "timelines": [
 {
 "tli": 1,
 "parent-tli": 0,
 "switchpoint": "0/0",
 "min-segno": "0000000010000000000000000001B",
 "max-segno": "00000000100000000000000000042",
 "n-segments": 40,
```

```

"size": 671088640,
"zratio": 1.00,
"closest-backup-id": "",
"status": "OK",
"lost-segments": [],
"backups": [
 {
 "id": "SBOL94",
 "status": "OK",
 "start-time": "2024-04-09 18:19:52+03",
 "backup-mode": "FULL",
 "wal": "ARCHIVE",
 "compress-alg": "lz4",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,
 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 2,
 "parent-tli": 0,
 "start-lsn": "0/41000028",
 "stop-lsn": "0/420000C0",
 "end-time": "2024-04-09 18:19:58+03",
 "end-validation-time": "2024-04-09 18:19:59+03",
 "recovery-xid": 757,
 "recovery-time": "2024-04-09 18:19:56.603355+03",
 "data-bytes": 395651278,
 "wal-bytes": 16777216,
 "uncompressed-bytes": 578552566,
 "pgdata-bytes": 578552248,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLsv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 3862224379
 },
 {
 "id": "SBOL8S",
 "status": "OK",
 "start-time": "2024-04-09 18:19:40+03",
 "parent-backup-id": "SBOL8G",
 "backup-mode": "DELTA",
 "wal": "STREAM",
 "compress-alg": "lz4",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,
 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 1,
 "parent-tli": 1,
 "start-lsn": "0/3C003020",
 "stop-lsn": "0/3E8D4930",
 "end-time": "2024-04-09 18:19:43+03",

```

```

 "end-validation-time": "2024-04-09 18:19:44+03",
 "recovery-xid": 757,
 "recovery-time": "2024-04-09 18:19:43.707720+03",
 "data-bytes": 119350434,
 "wal-bytes": 67108864,
 "uncompressed-bytes": 170044286,
 "pgdata-bytes": 578552248,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLsv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 1259851036
 },
 {
 "id": "SBOL8G",
 "status": "OK",
 "start-time": "2024-04-09 18:19:28+03",
 "parent-backup-id": "SBOL83",
 "backup-mode": "PTRACK",
 "wal": "STREAM",
 "compress-alg": "zlib",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,
 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 1,
 "parent-tli": 1,
 "start-lsn": "0/31000028",
 "stop-lsn": "0/310029E0",
 "end-time": "2024-04-09 18:19:32+03",
 "end-validation-time": "2024-04-09 18:19:33+03",
 "recovery-xid": 756,
 "recovery-time": "2024-04-09 18:19:32.594670+03",
 "data-bytes": 31218302,
 "wal-bytes": 16777216,
 "uncompressed-bytes": 69610366,
 "pgdata-bytes": 510263736,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLsv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 2293248595
 },
 {
 "id": "SBOL83",
 "status": "OK",
 "start-time": "2024-04-09 18:19:15+03",
 "parent-backup-id": "SBOL7P",
 "backup-mode": "PAGE",
 "wal": "STREAM",
 "compress-alg": "pglz",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,

```

```

 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 1,
 "parent-tli": 1,
 "start-lsn": "0/29000028",
 "stop-lsn": "0/2A0000F8",
 "end-time": "2024-04-09 18:19:22+03",
 "end-validation-time": "2024-04-09 18:19:22+03",
 "recovery-xid": 755,
 "recovery-time": "2024-04-09 18:19:22.269595+03",
 "data-bytes": 48394744,
 "wal-bytes": 33554432,
 "uncompressed-bytes": 69577598,
 "pgdata-bytes": 441975224,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLsv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 343227200
 },
 {
 "id": "SBOL7P",
 "status": "OK",
 "start-time": "2024-04-09 18:19:01+03",
 "backup-mode": "FULL",
 "wal": "STREAM",
 "compress-alg": "zstd",
 "compress-level": 1,
 "from-replica": "false",
 "block-size": 8192,
 "xlog-block-size": 8192,
 "checksum-version": 1,
 "program-version": "2.7.3",
 "server-version": "17",
 "current-tli": 1,
 "parent-tli": 0,
 "start-lsn": "0/22000028",
 "stop-lsn": "0/220001C8",
 "end-time": "2024-04-09 18:19:07+03",
 "end-validation-time": "2024-04-09 18:19:09+03",
 "recovery-xid": 754,
 "recovery-time": "2024-04-09 18:19:06.557301+03",
 "data-bytes": 151177272,
 "wal-bytes": 16777216,
 "uncompressed-bytes": 373695222,
 "pgdata-bytes": 373694904,
 "primary_conninfo": "user=backup channel_binding=prefer
host=localhost port=5432 sslmode=prefer sslcompression=0 sslcertmode=allow sslsni=1
ssl_min_protocol_version=TLsv1.2 gssencmode=prefer krbsrvname=postgres gssdelegation=0
target_session_attrs=any load_balance_hosts=disable",
 "content-crc": 1636300818
 }
]
}
]
}
]
]

```

Most fields are consistent with the plain format, with some exceptions:

- The size is in bytes.
- The `closest-backup-id` attribute contains the ID of the most recent valid backup that belongs to one of the previous timelines. You can use this backup to perform point-in-time recovery to this timeline. If such a backup does not exist, this string is empty.
- The `lost-segments` array provides with information about intervals of missing segments in `DEGRADED` timelines. In `OK` timelines, the `lost-segments` array is empty.
- The `backups` array lists all backups belonging to the timeline. If the timeline has no backups, this array is empty.

## Configuring Retention Policy

With `pg_probackup`, you can configure retention policy to remove redundant backups, clean up unneeded WAL files, as well as pin specific backups to ensure they are kept for the specified time, as explained in the sections below. All these actions can be combined together in any way.

### Removing Redundant Backups

By default, all backup copies created with `pg_probackup` are stored in the specified backup catalog. To save disk space, you can configure retention policy to remove redundant backup copies.

To configure retention policy, set one or more of the following variables in the `pg_probackup.conf` file via [set-config](#):

```
--retention-redundancy=redundancy
```

Specifies **the number of full backup copies** to keep in the backup catalog.

```
--retention-window=window
```

Defines the earliest point in time for which `pg_probackup` can complete the recovery. This option is set in **the number of days** from the current moment. For example, if `retention-window=7`, `pg_probackup` must keep at least one backup copy that is older than seven days, with all the corresponding WAL files, and all the backups that follow.

If both `--retention-redundancy` and `--retention-window` options are set, both these conditions have to be taken into account when purging the backup catalog. For example, if you set `--retention-redundancy=2` and `--retention-window=7`, `pg_probackup` has to keep two full backup copies, as well as all the backups required to ensure recoverability for the last seven days:

```
pg_probackup set-config -B backup_dir --instance=instance_name --retention-redundancy=2
--retention-window=7
```

It is recommended to always keep at least two last parent full backups to avoid errors when creating incremental backups.

To clean up the backup catalog in accordance with retention policy, you have to run the [delete](#) command with [retention flags](#), as shown below, or use the [backup](#) command with these flags to process the outdated backup copies right when the new backup is created.

For example, to remove all backup copies that no longer satisfy the defined retention policy, run the following command with the `--delete-expired` flag:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-expired
```

If you would like to also remove the WAL files that are no longer required for any of the backups, you should also specify the `--delete-wal` flag:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-expired --delete-wal
```

You can also set or override the current retention policy by specifying `--retention-redundancy` and `--retention-window` options directly when running `delete` or `backup` commands:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-expired --
retention-window=7 --retention-redundancy=2
```

Since incremental backups require that their parent full backup and all the preceding incremental backups are available, if any of such backups expire, they still cannot be removed while at least one incremental backup in this chain satisfies the retention policy. To avoid keeping expired backups that are still required to restore an active incremental one, you can merge them with this backup using the `--merge-expired` flag when running `backup` or `delete` commands.

Suppose you have backed up the `node` instance in the `backup_dir` directory, with the `--retention-window` option set to 6 and `--retention-redundancy` option set to 2, and you have the following backups available on August 01, 2024:

```
BACKUP INSTANCE 'node'
```

| Instance                   | Version | ID     | Recovery Time          | Mode       | WAL Mode | TLI | Time | Data   |
|----------------------------|---------|--------|------------------------|------------|----------|-----|------|--------|
| WAL                        | Zalg    | Zratio | Start LSN              | Stop LSN   | Status   |     |      |        |
| node                       | 17      | SHJ1N9 | 2024-08-01 09:50:00+03 | FULL       | ARCHIVE  | 1/0 | 5s   | 13MB   |
| 16MB                       | zstd    | 2,81   | 0/1D000028             | 0/1E0000C0 | OK       |     |      |        |
| node                       | 17      | SHJ1N8 | 2024-08-01 09:49:59+03 | DELTA      | ARCHIVE  | 1/1 | 5s   | 6432kB |
| 16MB                       | zstd    | 1,06   | 0/1A000028             | 0/1B0000C0 | OK       |     |      |        |
| node                       | 17      | SHH6Z8 | 2024-07-31 09:49:59+03 | PAGE       | ARCHIVE  | 1/1 | 5s   | 6431kB |
| 16MB                       | zstd    | 1,06   | 0/17000028             | 0/180000C0 | OK       |     |      |        |
| node                       | 17      | SHFCB6 | 2024-07-30 09:49:57+03 | FULL       | ARCHIVE  | 1/0 | 5s   | 12MB   |
| 16MB                       | zstd    | 2,83   | 0/14000028             | 0/150000C0 | OK       |     |      |        |
| node                       | 17      | SH9SB5 | 2024-07-27 09:49:56+03 | PAGE       | ARCHIVE  | 1/1 | 5s   | 6432kB |
| 16MB                       | zstd    | 1,06   | 0/11000028             | 0/120000C0 | OK       |     |      |        |
| -----retention window----- |         |        |                        |            |          |     |      |        |
| node                       | 17      | SH62Z5 | 2024-07-25 09:49:56+03 | DELTA      | ARCHIVE  | 1/1 | 5s   | 6431kB |
| 16MB                       | zstd    | 1,06   | 0/E000028              | 0/F0000C0  | OK       |     |      |        |
| node                       | 17      | SH48B3 | 2024-07-24 09:49:54+03 | FULL       | ARCHIVE  | 1/0 | 5s   | 12MB   |
| 16MB                       | zstd    | 2,86   | 0/B000028              | 0/C0000C0  | OK       |     |      |        |
| node                       | 17      | SGWTN3 | 2024-07-20 09:49:54+03 | PAGE       | ARCHIVE  | 1/1 | 5s   | 6432kB |
| 16MB                       | zstd    | 1,06   | 0/8000028              | 0/90000C0  | OK       |     |      |        |
| node                       | 17      | SGUYZ2 | 2024-07-19 09:49:53+03 | DELTA      | ARCHIVE  | 1/1 | 5s   | 6442kB |
| 16MB                       | zstd    | 1,07   | 0/5000028              | 0/60000C0  | OK       |     |      |        |
| node                       | 17      | SGT4B1 | 2024-07-18 09:49:52+03 | FULL       | ARCHIVE  | 1/0 | 5s   | 11MB   |
| 16MB                       | zstd    | 2,89   | 0/2000028              | 0/3003A28  | OK       |     |      |        |

If you run the `delete` command with the `--delete-expired` flag, the backups with IDs `SGT4B1`, `SGUYZ2`, and `SGWTN3` will be removed as they are expired both according to the retention window and due to redundancy (the required set of full backups has already been retained). `SGUYZ2` and `SGWTN3` will also be removed since the base full backup is expired.

Running the `delete` command with the `--merge-expired` flag will merge backups `SH48B3` and `SH62Z5` with `SH9SB5`. The merge will occur with `SH9SB5` as it is the first non-expired delta backup, which can be merged with expired delta backups `SH62Z5` and expired full backup `SH48B3`.

```
pg_probackup delete -B backup_dir --instance=node --delete-expired --merge-expired
pg_probackup show -B backup_dir
```

```
BACKUP INSTANCE 'node'
```

| Instance | Version | ID     | Recovery Time          | Mode       | WAL Mode | TLI | Time | Data |
|----------|---------|--------|------------------------|------------|----------|-----|------|------|
| WAL      | Zalg    | Zratio | Start LSN              | Stop LSN   | Status   |     |      |      |
| node     | 17      | SHJ1N9 | 2024-08-01 09:50:00+03 | FULL       | ARCHIVE  | 1/0 | 5s   | 13MB |
| 16MB     | zstd    | 2,81   | 0/1D000028             | 0/1E0000C0 | OK       |     |      |      |

|      |      |        |                        |       |         |     |    |        |
|------|------|--------|------------------------|-------|---------|-----|----|--------|
| node | 17   | SHJ1N8 | 2024-08-01 09:49:59+03 | DELTA | ARCHIVE | 1/1 | 5s | 6432kB |
| 16MB | zstd | 1,06   | 0/1A000028 0/1B0000C0  | OK    |         |     |    |        |
| node | 17   | SHH6Z8 | 2024-07-31 09:49:59+03 | PAGE  | ARCHIVE | 1/1 | 5s | 6431kB |
| 16MB | zstd | 1,06   | 0/17000028 0/180000C0  | OK    |         |     |    |        |
| node | 17   | SHFCB6 | 2024-07-30 09:49:57+03 | FULL  | ARCHIVE | 1/0 | 5s | 12MB   |
| 16MB | zstd | 2,83   | 0/14000028 0/150000C0  | OK    |         |     |    |        |
| node | 17   | SH9SB5 | 2024-07-27 09:49:56+03 | FULL  | ARCHIVE | 1/0 | 1s | 12MB   |
| 16MB | zstd | 2,84   | 0/11000028 0/120000C0  | OK    |         |     |    |        |

The `Time` field for the merged backup displays the time required for the merge.

### Pinning Backups

If you need to keep certain backups longer than the established retention policy allows, you can pin them for arbitrary time. For example:

```
pg_probackup set-backup -B backup_dir --instance=instance_name -i backup_id --ttl=30d
```

This command sets the expiration time of the specified backup to 30 days starting from the time indicated in its `recovery-time` attribute.

You can also explicitly set the expiration time for a backup using the `--expire-time` option. For example:

```
pg_probackup set-backup -B backup_dir --instance=instance_name -i backup_id --expire-time="2027-04-09 18:21:32+00"
```

Alternatively, you can use the `--ttl` and `--expire-time` options with the `backup` command to pin the newly created backup:

```
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --ttl=30d
pg_probackup backup -B backup_dir --instance=instance_name -b FULL --expire-time="2027-04-09 18:21:32+00"
```

To check if the backup is pinned, run the `show` command:

```
pg_probackup show -B backup_dir --instance=instance_name -i backup_id
```

If the backup is pinned, it has the `expire-time` attribute that displays its expiration time:

```
...
recovery-time = '2024-04-09 18:21:32+00'
expire-time = '2027-04-09 18:21:32+00'
data-bytes = 22288792
...
```

You can unpin the backup by setting the `--ttl` option to zero:

```
pg_probackup set-backup -B backup_dir --instance=instance_name -i backup_id --ttl=0
```

### Note

A pinned incremental backup implicitly pins all its parent backups. If you unpin such a backup later, its implicitly pinned parents will also be automatically unpinned.

### Configuring WAL Archive Retention Policy

When [continuous WAL archiving](#) is enabled, archived WAL segments can take a lot of disk space. Even if you delete old backup copies from time to time, the `--delete-wal` flag can purge only those WAL segments that do not apply to any of the remaining backups in the backup catalog. However, if point-in-time recovery is critical only for the most recent backups, you can configure WAL archive retention policy to keep WAL archive of limited depth and win back some more disk space.

To configure WAL archive retention policy, you have to run the `set-config` command with the `--wal-depth` option that specifies the number of backups that can be used for PITR. This setting applies to all



the timelines, so you should be able to perform PITR for the same number of backups on each timeline, if available. [Pinned backups](#) are not included into this count: if one of the latest backups is pinned, `pg_probackup` ensures that PITR is possible for one extra backup.

To remove WAL segments that do not satisfy the defined WAL archive retention policy, you simply have to run the [delete](#) or [backup](#) command with the `--delete-wal` flag. For archive backups, WAL segments between `Start LSN` and `Stop LSN` are always kept intact, so such backups remain valid regardless of the `--wal-depth` setting and can still be restored, if required.

You can also use the `--wal-depth` option with the [delete](#) and [backup](#) commands to override the previously defined WAL archive retention policy and purge old WAL segments on the fly.

Suppose you have backed up the `node` instance in the `backup_dir` directory and configured [continuous WAL archiving](#):

```
pg_probackup show -B backup_dir --instance=node
```

| Instance | Version | ID     | Recovery Time                 | Mode       | WAL Mode   | TLI    | Time |
|----------|---------|--------|-------------------------------|------------|------------|--------|------|
| Data     | WAL     | Zalg   | Zratio                        | Start LSN  | Stop LSN   | Status |      |
| node     | 17      | SBOLDA | 2024-04-09 18:22:23.147138+03 | DELTA      | STREAM     | 1/1    | 1s   |
| 1165kB   | 16MB    | zstd   | 1.09                          | 0/6F000028 | 0/6F000190 | OK     |      |
| node     | 17      | SBOLCY | 2024-04-09 18:22:16.079841+03 | FULL       | STREAM     | 1/0    | 10s  |
| 278MB    | 16MB    | zstd   | 2.46                          | 0/6D000028 | 0/6D000190 | OK     |      |
| node     | 17      | SBOLCW | 2024-04-09 18:22:10.154022+03 | DELTA      | STREAM     | 1/1    | 2s   |
| 1364kB   | 16MB    | zstd   | 1.01                          | 0/6B000028 | 0/6B000190 | OK     |      |
| node     | 17      | SBOLCS | 2024-04-09 18:22:07.521646+03 | DELTA      | STREAM     | 1/1    | 4s   |
| 78MB     | 16MB    | zstd   | 2.41                          | 0/69000028 | 0/69000190 | OK     |      |
| node     | 17      | SBOLCC | 2024-04-09 18:21:55.830115+03 | FULL       | STREAM     | 1/0    | 15s  |
| 278MB    | 96MB    | zstd   | 2.46                          | 0/600060C8 | 0/64FE6640 | OK     |      |
| node     | 17      | SBOLBW | 2024-04-09 18:21:38.399702+03 | FULL       | STREAM     | 1/0    | 12s  |
| 278MB    | 96MB    | zstd   | 2.46                          | 0/54001830 | 0/589E5908 | OK     |      |

You can check the state of the WAL archive by running the [show](#) command with the `--archive` flag:

```
pg_probackup show -B backup_dir --instance=node --archive
```

```
INFO: checking WAL file name "000000010000000000000000000052"
INFO: checking WAL file name "000000010000000000000000000053"
INFO: checking WAL file name "000000010000000000000000000054"
INFO: checking WAL file name "000000010000000000000000000054.00001830.backup"
INFO: checking WAL file name "000000010000000000000000000055"
INFO: checking WAL file name "000000010000000000000000000056"
INFO: checking WAL file name "000000010000000000000000000057"
INFO: checking WAL file name "000000010000000000000000000058"
INFO: checking WAL file name "000000010000000000000000000059"
INFO: checking WAL file name "00000001000000000000000000005A"
INFO: checking WAL file name "00000001000000000000000000005B"
INFO: checking WAL file name "00000001000000000000000000005C"
INFO: checking WAL file name "00000001000000000000000000005D"
INFO: checking WAL file name "00000001000000000000000000005E"
INFO: checking WAL file name "00000001000000000000000000005F"
INFO: checking WAL file name "000000010000000000000000000060"
INFO: checking WAL file name "000000010000000000000000000060.000060C8.backup"
INFO: checking WAL file name "000000010000000000000000000061"
INFO: checking WAL file name "000000010000000000000000000062"
INFO: checking WAL file name "000000010000000000000000000063"
INFO: checking WAL file name "000000010000000000000000000064"
INFO: checking WAL file name "000000010000000000000000000065"
```

## Additional Supplied Programs

```
INFO: checking WAL file name "00000001000000000000000066"
INFO: checking WAL file name "00000001000000000000000067"
INFO: checking WAL file name "00000001000000000000000068"
INFO: checking WAL file name "00000001000000000000000069"
INFO: checking WAL file name "00000001000000000000000069.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006A"
INFO: checking WAL file name "0000000100000000000000006B"
INFO: checking WAL file name "0000000100000000000000006B.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006C"
INFO: checking WAL file name "0000000100000000000000006D"
INFO: checking WAL file name "0000000100000000000000006D.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006E"
INFO: checking WAL file name "0000000100000000000000006F"
INFO: checking WAL file name "0000000100000000000000006F.00000028.backup"
```

ARCHIVE INSTANCE 'node'

| TLI      | Parent TLI | Switchpoint | Min Segno                  | Max Segno                  | N  |
|----------|------------|-------------|----------------------------|----------------------------|----|
| segments | Size       | Zratio      | N backups                  | Status                     |    |
| 1        | 0          | 0/0         | 00000001000000000000000052 | 0000000100000000000000006F | 30 |
|          | 480MB      | 1.00        | 6                          | OK                         |    |

**WAL purge without --wal-depth cannot achieve much, only one segment is removed:**

```
pg_probackup delete -B backup_dir --instance=node --delete-wal
```

```
INFO: checking WAL file name "00000001000000000000000054"
INFO: checking WAL file name "00000001000000000000000054.00001830.backup"
INFO: checking WAL file name "00000001000000000000000055"
INFO: checking WAL file name "00000001000000000000000056"
INFO: checking WAL file name "00000001000000000000000057"
INFO: checking WAL file name "00000001000000000000000058"
INFO: checking WAL file name "00000001000000000000000059"
INFO: checking WAL file name "0000000100000000000000005A"
INFO: checking WAL file name "0000000100000000000000005B"
INFO: checking WAL file name "0000000100000000000000005C"
INFO: checking WAL file name "0000000100000000000000005D"
INFO: checking WAL file name "0000000100000000000000005E"
INFO: checking WAL file name "0000000100000000000000005F"
INFO: checking WAL file name "00000001000000000000000060"
INFO: checking WAL file name "00000001000000000000000060.000060C8.backup"
INFO: checking WAL file name "00000001000000000000000061"
INFO: checking WAL file name "00000001000000000000000062"
INFO: checking WAL file name "00000001000000000000000063"
INFO: checking WAL file name "00000001000000000000000064"
INFO: checking WAL file name "00000001000000000000000065"
INFO: checking WAL file name "00000001000000000000000066"
INFO: checking WAL file name "00000001000000000000000067"
INFO: checking WAL file name "00000001000000000000000068"
INFO: checking WAL file name "00000001000000000000000069"
INFO: checking WAL file name "00000001000000000000000069.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006A"
INFO: checking WAL file name "0000000100000000000000006B"
INFO: checking WAL file name "0000000100000000000000006B.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006C"
INFO: checking WAL file name "0000000100000000000000006D"
INFO: checking WAL file name "0000000100000000000000006D.00000028.backup"
INFO: checking WAL file name "0000000100000000000000006E"
```

```
INFO: checking WAL file name "00000001000000000000000006F"
INFO: checking WAL file name "00000001000000000000000006F.00000028.backup"
```

```
ARCHIVE INSTANCE 'node'
```

| TLI segments | Parent TLI Size | Switchpoint Zratio | Min Segno N backups | Segno Status                    | Max Segno                 | N  |
|--------------|-----------------|--------------------|---------------------|---------------------------------|---------------------------|----|
| 1            | 0               | 0/0                | 6                   | 0000000100000000000000054<br>OK | 000000010000000000000006F | 28 |

If you would like, for example, to keep only those WAL segments that can be applied to the latest valid backup, set the `--wal-depth` option to 1:

```
pg_probackup delete -B backup_dir --instance=node --delete-wal --wal-depth=1
```

```
INFO: checking WAL file name "00000001000000000000000006F"
INFO: checking WAL file name "00000001000000000000000006F.00000028.backup"
```

```
ARCHIVE INSTANCE 'node'
```

| TLI segments | Parent TLI Size | Switchpoint Zratio | Min Segno N backups | Segno Status                    | Max Segno                 | N |
|--------------|-----------------|--------------------|---------------------|---------------------------------|---------------------------|---|
| 1            | 0               | 0/0                | 6                   | 000000010000000000000006F<br>OK | 000000010000000000000006F | 1 |

Alternatively, you can use the `--wal-depth` option with the `backup` command:

```
pg_probackup backup -B backup_dir --instance=node -b DELTA --wal-depth=1 --delete-wal
```

```
INFO: checking WAL file name "000000010000000000000000071"
INFO: checking WAL file name "000000010000000000000000071.00000028.backup"
```

```
ARCHIVE INSTANCE 'node'
```

| TLI segments | Parent TLI Size | Switchpoint Zratio | Min Segno N backups | Segno Status                    | Max Segno                 | N |
|--------------|-----------------|--------------------|---------------------|---------------------------------|---------------------------|---|
| 1            | 0               | 0/0                | 7                   | 0000000100000000000000071<br>OK | 0000000100000000000000071 | 1 |

## Merging Backups

As you take more and more incremental backups, the total size of the backup catalog can substantially grow. To save disk space, you can merge incremental backups to their parent full backup by running the `merge` command, specifying the backup ID of the most recent incremental backup you would like to merge:

```
pg_probackup merge -B backup_dir --instance=instance_name -i backup_id
```

This command merges backups that belong to a common incremental backup chain. If you specify a full backup, it will be merged with its first incremental backup. If you specify an incremental backup, it will be merged to its parent full backup, together with all incremental backups between them. Once the merge is complete, the full backup takes in all the merged data, and the incremental backups are removed as redundant. Thus, the merge operation is virtually equivalent to retaking a full backup and removing all the outdated backups, but it allows you to save much time, especially for large data volumes, as well as I/O and network traffic if you are using `pg_probackup` in the [remote](#) mode.

Before the merge, `pg_probackup` validates all the affected backups to ensure that they are valid. You can check the current backup status by running the `show` command with the backup ID:

```
pg_probackup show -B backup_dir --instance=instance_name -i backup_id
```

If the merge is still in progress, the backup status is displayed as `MERGING`. For full backups, it can also be shown as `MERGED` while the metadata is being updated at the final stage of the merge. The merge is idempotent, so you can restart the merge if it was interrupted.

## Deleting Backups

To delete a backup that is no longer required, run the following command:

```
pg_probackup delete -B backup_dir --instance=instance_name -i backup_id
```

This command will delete the backup with the specified `backup_id`, together with all the incremental backups that descend from `backup_id`, if any. This way you can delete some recent incremental backups, retaining the underlying full backup and some of the incremental backups that follow it.

To delete obsolete WAL files that are not necessary to restore any of the remaining backups, use the `--delete-wal` flag:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-wal
```

To delete backups that are expired according to the current retention policy, use the `--delete-expired` flag:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-expired
```

Expired backups cannot be removed while at least one incremental backup that satisfies the retention policy is based on them. If you would like to minimize the number of backups still required to keep incremental backups valid, specify the `--merge-expired` flag when running this command:

```
pg_probackup delete -B backup_dir --instance=instance_name --delete-expired --merge-expired
```

In this case, `pg_probackup` searches for the oldest incremental backup that satisfies the retention policy and merges this backup with the underlying full and incremental backups that have already expired, thus making it a full backup. Once the merge is complete, the remaining expired backups are deleted.

Before merging or deleting backups, you can run the `delete` command with the `--dry-run` flag, which displays the status of all the available backups according to the current retention policy, without performing any irreversible actions.

To delete all backups with specific status, use the `--status`:

```
pg_probackup delete -B backup_dir --instance=instance_name --status=ERROR
```

Deleting backups by status ignores established retention policies.

## Cloning and Synchronizing Postgres Pro Instance

`pg_probackup` can create a copy of a Postgres Pro instance directly, without using the backup catalog. To do this, you can run the `catchup` command. It can be useful in the following cases:

- To add a new standby server.

If the data directory of the destination instance is empty, the `catchup` command works similarly, but it can be faster if run in parallel mode.

- To have a fallen-behind standby server “catch up” with the primary.

Under write-intensive load, replicas may fail to replay WAL fast enough to keep up with the primary and hence may lag behind. A usual solution to create a new replica and switch to it requires a lot of extra space and data transfer. The `catchup` command allows you to update an existing replica much faster by bringing differences from the primary.

`catchup` is different from other `pg_probackup` operations:

- The backup catalog is not required.

- STREAM WAL delivery mode is only supported.
- Copying external directories is not supported.
- DDL commands `CREATE TABLESPACE/DROP TABLESPACE` cannot be run simultaneously with `catchup`.
- `catchup` takes configuration files, such as `postgresql.conf`, `postgresql.auto.conf`, or `pg_hba.conf`, from the source server and overwrites them on the target server. The `--exclude-path` option allows you to keep the configuration files intact.

To prepare for cloning/synchronizing a Postgres Pro instance, set up the source server as follows:

- [Configure the database cluster](#) for the instance to copy.
- To copy from a remote server, [configure the remote mode](#).
- To use the PTRACK `catchup` mode, [set up PTRACK backups](#).

Before cloning/synchronizing a Postgres Pro instance, ensure that the source server is running and accepting connections. To clone/sync a Postgres Pro instance, on the server with the destination instance, you can run the `catchup` command as follows:

```
pg_probackup catchup -b catchup_mode --source-pgdata=path_to_pgdata_on_remote_server --
destination-pgdata=path_to_local_dir --stream [connection_options] [remote_options]
```

Where `catchup_mode` can take one of the following values:

- `FULL` — creates a full copy of the Postgres Pro instance. The data directory of the destination instance must be empty for this mode.
- `DELTA` — reads all data files in the data directory and creates an incremental copy for pages that have changed since the destination instance was shut down.
- `PTRACK` — tracking page changes on the fly, only reads and copies pages that have changed since the point of divergence of the source and destination instances.

### Warning

PTRACK `catchup` mode requires PTRACK not earlier than 2.0 and hence, Postgres Pro not earlier than 11.

By specifying the `--stream` option, you can set [STREAM](#) WAL delivery mode of copying, which will include all the necessary WAL files by streaming them from the server via replication protocol.

You can use [connection options](#) to specify the connection to the source database cluster. If it is located on a different server, also specify [remote options](#).

If the source database cluster contains tablespaces that must be located in a different directory, additionally specify the `--tablespace-mapping` option:

```
pg_probackup catchup -b catchup_mode --source-pgdata=path_to_pgdata_on_remote_server --
destination-pgdata=path_to_local_dir --stream --tablespace-mapping=OLDDIR=NEWDIR
```

To run the `catchup` command on parallel threads, specify the number of threads with the `--threads` option:

```
pg_probackup catchup -b catchup_mode --source-pgdata=path_to_pgdata_on_remote_server --
destination-pgdata=path_to_local_dir --stream --threads=num_threads
```

Before cloning/synchronising a Postgres Pro instance, you can run the `catchup` command with the `--dry-run` flag to estimate the size of data files to be transferred, but make no changes on disk:

```
pg_probackup catchup -b catchup_mode --source-pgdata=path_to_pgdata_on_remote_server --
destination-pgdata=path_to_local_dir --stream --dry-run
```

For example, assume that a remote standby server with the Postgres Pro instance having `/replica-pg-data` data directory has fallen behind. To sync this instance with the one in `/master-pgdata` data directory, you can run the `catchup` command in the `PTRACK` mode on four parallel threads as follows:

```
pg_probackup catchup --source-pgdata=/master-pgdata --destination-pgdata=/replica-pgdata -p 5432 -d postgres -U remote-postgres-user --stream --backup-mode=PTRACK --remote-host=remote-hostname --remote-user=remote-unix-username -j 4 --exclude-path=postgresql.conf --exclude-path=postgresql.auto.conf --exclude-path=pg_hba.conf --exclude-path=pg_ident.conf
```

Note that in this example, the configuration files will not be overwritten during synchronization.

Another example shows how you can add a new remote standby server with the Postgres Pro data directory `/replica-pgdata` by running the `catchup` command in the `FULL` mode on four parallel threads:

```
pg_probackup catchup --source-pgdata=/master-pgdata --destination-pgdata=/replica-pgdata -p 5432 -d postgres -U remote-postgres-user --stream --backup-mode=FULL --remote-host=remote-hostname --remote-user=remote-unix-username -j 4
```

## Command-Line Reference

### Commands

This section describes `pg_probackup` commands. Optional parameters are enclosed in square brackets. For detailed parameter descriptions, see the section [Options](#).

#### version

```
pg_probackup version [--format=json]
```

Prints `pg_probackup` version and edition, as well as Postgres Pro version and edition.

If `--format=json` is specified, the output is printed in the JSON format. This may be needed for native integration with JSON-based applications, such as PPEM. Example of a JSON output:

```
pg_probackup version --format=json
{
 "pg_probackup": {
 "version": "2.8.2",
 "edition": "enterprise"
 },
 "database": {
 "type": "Postgres Pro Enterprise",
 "version": "16.3.1"
 },
 "compressions": [zlib, pglz, lz4, zstd]
}
```

#### help

```
pg_probackup help [command]
```

Displays the synopsis of `pg_probackup` commands. If one of the `pg_probackup` commands is specified, shows detailed information about the options that can be used with this command.

#### init

```
pg_probackup init -B backup_dir [--skip-if-exists] [s3_options] [--help]
[logging_options]
```

Initializes the backup catalog in `backup_dir` that will store backup copies, WAL archive, and meta information for the backed up database clusters. If the specified `backup_dir` already exists, it must be

empty. Otherwise, `pg_probackup` displays a corresponding error message. You can ignore this error by specifying the `--skip-if-exists` option. Although the backup will not be initialized, the application will return 0 code.

For details, see the section [Initializing the Backup Catalog](#).

#### **add-instance**

```
pg_probackup add-instance -B backup_dir -D data_dir --instance=instance_name
[--skip-if-exists] [s3_options] [--help] [logging_options]
```

Initializes a new backup instance inside the backup catalog `backup_dir` and generates the `pg_probackup.conf` configuration file that controls `pg_probackup` settings for the cluster with the specified `data_dir` data directory. If the catalog was already initialized, you can ignore the error by specifying `--skip-if-exists`.

For details, see the section [Adding a New Backup Instance](#).

#### **del-instance**

```
pg_probackup del-instance -B backup_dir --instance=instance_name [s3_options] [--help]
[logging_options]
```

Deletes all backups and WAL files associated with the specified instance.

#### **set-config**

```
pg_probackup set-config -B backup_dir --instance=instance_name
[--help] [--pgdata=pgdata-path]
[--retention-redundancy=redundancy] [--retention-window=window] [--wal-depth=wal_depth]
[--compress-algorithm=compression_algorithm] [--compress-level=compression_level]
[-d dbname] [-h host] [-p port] [-U username]
[--archive-timeout=timeout] [--external-dirs=external_directory_path]
[--restore-command=cmdline]
[remote_options] [remote_wal_archive_options] [logging_options] [s3_options]
```

Adds the specified connection, compression, retention, logging, and external directory settings into the `pg_probackup.conf` configuration file, or modifies the previously defined values.

For all available settings, see the [Options](#) section.

It is **not recommended** to edit `pg_probackup.conf` manually.

#### **set-backup**

```
pg_probackup set-backup -B backup_dir --instance=instance_name -i backup_id
{--ttl=ttl | --expire-time=time}
[--note=backup_note] [s3_options] [--help] [logging_options]
```

Sets the provided backup-specific settings into the `backup.control` configuration file, or modifies the previously defined values.

`--note=backup_note`

Sets the text note for backup copy. If `backup_note` contains newline characters, then only the substring before the first newline character will be saved. The maximum size of a text note is 1 KB. The `'none'` value removes the current note.

For all available pinning settings, see the section [Pinning Options](#).

#### **show-config**

```
pg_probackup show-config -B backup_dir --instance instance_name [--format=plain|json]
[s3_options]
[--no-scale-units] [logging_options]
```

Displays all the current `pg_probackup` configuration settings, including those that are specified in the `pg_probackup.conf` configuration file located in the `backup_dir/backups/instance_name` directory and those that were provided on a command line. You can specify the `--format=json` option to get the result in the JSON format. By default, configuration settings are shown as plain text.

You can also specify the `--no-scale-units` option to display time and memory configuration settings in their base (unscaled) units. Otherwise, the values are scaled to larger units for optimal display. For example, if `archive-timeout` is 300, then 5min is displayed, but if `archive-timeout` is 301, then 301s is displayed. Also, if the `--no-scale-units` option is specified, configuration settings are displayed without units and for the JSON format, numeric and boolean values are not enclosed in quotes. This facilitates parsing the output.

To edit `pg_probackup.conf`, use the [set-config](#) command.

#### show

```
pg_probackup show -B backup_dir
[--help] [--instance=instance_name [-i backup_id | --archive]] [--format=plain|json]
[--no-color] [--show-symlinks] [s3_options]
[logging_options]
```

Shows the contents of the backup catalog. If `instance_name` and `backup_id` are specified, shows detailed information about this backup. If the `--archive` option is specified, shows the contents of WAL archive of the backup catalog.

By default, the contents of the backup catalog is shown as plain text. You can specify the `--format=json` option to get the result in the JSON format. If `--no-color` flag is used, then the output is not colored.

If the `--show-symlinks` option is specified, the command also shows the links between [merged](#) backups and the original full backups that incremental backups were merged to.

For details on usage, see the sections [Managing the Backup Catalog](#) and [Viewing WAL Archive Information](#).

#### backup

```
pg_probackup backup -B backup_dir -b backup_mode --instance=instance_name
[--help] [-j num_threads] [--progress]
[--backup-threads num_threads] [--validate-threads num_threads]
[-C] [--stream [-S slot_name] [--temp-slot[=true|false|on|off]]] [--backup-pg-log]
[--no-validate] [--skip-block-validation]
[-w --no-password] [-W --password]
[--write-rate-limit=bitrate]
[--archive-timeout=timeout] [--external-dirs=external_directory_path]
[--no-sync] [--note=backup_note]
[connection_options] [compression_options] [remote_options]
[retention_options] [pinning_options] [logging_options] [s3_options]
```

Creates a backup copy of the Postgres Pro instance.

`-b mode`  
`--backup-mode=mode`

Specifies the backup mode to use. Possible values are: [FULL](#), [DELTA](#), [PAGE](#), and [PTRACK](#).

`--backup-threads num_threads`

Specifies the number of threads for copying files. Overrides the `j/--threads` option for file copying.

`--validate-threads num_threads`

Specifies the number of threads for the backup validation. Overrides the `j/--threads` option for the backup validation.



`-C`  
`--smooth-checkpoint`  
Spreads out the checkpoint over a period of time. By default, `pg_probackup` tries to complete the checkpoint as soon as possible.

`--stream`  
Makes a **STREAM** backup, which includes all the necessary WAL files by streaming them from the database server via replication protocol.

`--temp-slot [=true|false|on|off]`  
Creates a *temporary* physical replication slot for streaming WAL from the backed up Postgres Pro instance. `--temp-slot` is enabled by default. It ensures that all the required WAL segments remain available if WAL is rotated while the backup is in progress. This flag can only be used together with the `--stream` flag. The default slot name is `pg_probackup_slot`. To change it, use the `--slot/-S` option and explicitly specify `--temp-slot` or `--temp-slot=true|on`.

`-S slot_name`  
`--slot=slot_name`  
Specifies the replication slot to connect to for WAL streaming. This option can only be used together with the `--stream` flag.

`--backup-pg-log`  
Includes the log directory into the backup. This directory usually contains log messages. By default, log directory is excluded.

`-E external_directory_path`  
`--external-dirs=external_directory_path`  
Includes the specified directory into the backup by recursively copying its contents into a separate subdirectory in the backup catalog. This option is useful to back up scripts, SQL dump files, and configuration files located outside of the data directory. If you would like to back up several external directories, separate their paths by a colon on Unix and a semicolon on Windows.

`--write-rate-limit=bitrate`  
Sets the rate of writing data to disk, in MBps or GBps. The default unit is MBps. For example: `--write-rate-limit=1GBps` or `--write-rate-limit=100` (MBps). The default value is 0 — no limitation.

If this option is specified, the following information is displayed at the end of the backup:

- `written` — the amount of data written, in MB.
- `total time` — the time that elapsed between the first and last writes, in seconds. Note that this is not the total backup time.
- `sleep time` — the amount of forced delay time, in seconds.
- `average rate` — the actual average write rate, in MBps.

For example:

```
INFO: Rate limit: written 14975.445 MB, total time 17.163 s, sleep time 2.370 s,
 average rate 872.560715 MBps
```

`--archive-timeout=wait_time`  
Sets the timeout for WAL segment archiving and streaming, in seconds. By default, `pg_probackup` waits 300 seconds.

`--skip-block-validation`  
Disables block-level checksum verification to speed up the backup process.

`--no-validate`

Skips automatic validation after the backup is taken. You can use this flag if you validate backups regularly and would like to save time when running backup operations.

It is recommended to use this flag when creating a backup to an S3 storage. Due to some features of S3 storages, automatic validation may appear incorrect in this case. Skip automatic validation and then perform validation using a separate [validate](#) command.

`--no-sync`

Do not sync backed up files to disk. You can use this flag to speed up the backup process. Using this flag can result in data corruption in case of operating system or hardware crash. If you use this option, it is recommended to run the [validate](#) command once the backup is complete to detect possible issues.

`--note=backup_note`

Sets the text note for backup copy. If *backup\_note* contain newline characters, then only substring before first newline character will be saved. Max size of text note is 1 KB. The *'none'* value removes current note.

For more details of the command settings, see sections [Common Options](#), [Connection Options](#), [Retention Options](#), [Pinning Options](#), [Remote Mode Options](#), [Compression Options](#), [Logging Options](#), and [S3 Options](#).

For details on usage, see the section [Creating a Backup](#).

## restore

```
pg_probackup restore -B backup_dir --instance=instance_name
[--help] [--dry-run] [-D data_dir] [-i backup_id]
[-j num_threads] [--progress]
[-T OLDDIR=NEWDIR] [--external-mapping=OLDDIR=NEWDIR] [--skip-external-dirs]
[-R | --restore-as-replica] [--no-validate] [--skip-block-validation]
[--force] [--no-sync]
[--restore-command=cmdline]
[--primary-conninfo=primary_conninfo]
[-S | --primary-slot-name=slot_name]
[-X wal_dir | --waldir=wal_dir]
[recovery_target_options] [logging_options] [remote_options]
[partial_restore_options] [remote_wal_archive_options] [s3_options]
```

Restores the Postgres Pro instance from a backup copy located in the *backup\_dir* backup catalog. If you specify a [recovery target option](#), `pg_probackup` finds the closest backup and restores it to the specified recovery target. If neither the backup ID nor recovery target options are provided, `pg_probackup` uses the most recent backup to perform the recovery.

`-R`

`--restore-as-replica`

Creates a minimal recovery configuration file to facilitate setting up a standby server. If the replication connection requires a password, you must specify the password manually in the [primary\\_conninfo](#) parameter as it is not included. For Postgres Pro 11 or lower, recovery settings are written into the `recovery.conf` file. Starting from Postgres Pro 12, `pg_probackup` writes these settings into the `probackup_recovery.conf` file in the data directory, and then includes them into the `postgresql.auto.conf` when the cluster is started.

`--primary-conninfo=primary_conninfo`

Sets the [primary\\_conninfo](#) parameter to the specified value. This option will be ignored unless the `-R` flag is specified.

Example: `--primary-conninfo="host=192.168.1.50 port=5432 user=foo password=foopass"`

`-S`  
`--primary-slot-name=slot_name`  
Sets the [primary\\_slot\\_name](#) parameter to the specified value. This option will be ignored unless the `-R` flag is specified.

`-T OLDDIR=NEWDIR`  
`--tablespace-mapping=OLDDIR=NEWDIR`  
Relocates the tablespace from the `OLDDIR` to the `NEWDIR` directory at the time of recovery. Both `OLDDIR` and `NEWDIR` must be absolute paths. If the path contains the equals sign (=), escape it with a backslash. This option can be specified multiple times for multiple tablespaces.

`--external-mapping=OLDDIR=NEWDIR`  
Relocates an external directory included into the backup from the `OLDDIR` to the `NEWDIR` directory at the time of recovery. Both `OLDDIR` and `NEWDIR` must be absolute paths. If the path contains the equals sign (=), escape it with a backslash. This option can be specified multiple times for multiple directories.

`--skip-external-dirs`  
Skip external directories included into the backup with the `--external-dirs` option. The contents of these directories will not be restored.

`--skip-block-validation`  
Disables block-level checksum verification to speed up validation. During automatic validation before the restore only file-level checksums will be verified.

`--no-validate`  
Skips backup validation. You can use this flag if you validate backups regularly and would like to save time when running restore operations.

`--restore-command=cmdline`  
Sets the [restore\\_command](#) parameter to the specified command. For example: `--restore-command='cp /mnt/server/archivedir/%f "%p"'`

`--force`  
Allows to ignore an invalid status of the backup. You can use this flag if you need to restore the Postgres Pro cluster from a corrupt or an invalid backup. Use with caution. If `PGDATA` contains a non-empty directory with system ID different from that of the backup being restored, [incremental restore](#) with this flag overwrites the directory contents (while an error occurs without the flag). If tablespaces are remapped through the `--tablespace-mapping` option into non-empty directories, the contents of such directories will be deleted.

`--no-sync`  
Do not sync restored files to disk. You can use this flag to speed up restore process. Using this flag can result in data corruption in case of operating system or hardware crash. If it happens, you have to run the [restore](#) command again.

`-X wal_dir`  
`--waldir=wal_dir`  
Sets the directory to write WAL files to. By default WAL files will be placed in the `pg_wal` subdirectory of the target directory, but this option can be used to place them elsewhere. `wal_dir` must be an absolute path, which must not already exist, but if it does, it must be empty.

For more details of the command settings, see sections [Common Options](#), [Recovery Target Options](#), [Remote Mode Options](#), [Remote WAL Archive Options](#), [Logging Options](#), [Partial Restore Options](#), and [S3 Options](#).

For details on usage, see the section [Restoring a Cluster](#).

### checkdb

```
pg_probackup checkdb
[-B backup_dir] [--instance=instance_name] [-D data_dir]
[--help] [-j num_threads] [--progress]
[--amcheck [--skip-block-validation] [--checkunique] [--heapallindexed]]
[connection_options] [logging_options]
[s3_options]
```

Verifies the Postgres Pro database cluster correctness by detecting physical and logical corruption.

For the command to work correctly, when the backup instance is created in the S3 storage, you must specify [S3 options](#) on the command line or through environment variables.

`--amcheck`

Performs logical verification of indexes for the specified Postgres Pro instance if no corruption was found while checking data files. You must have the [amcheck](#) extension or the `amcheck_next` extension installed in the database to check its indexes. For databases without `amcheck`, index verification will be skipped. Additional options `--checkunique` and `--heapallindexed` are effective depending on the version of `amcheck` installed.

`--checkunique`

Verifies unique constraints during logical verification of indexes. You can use this flag only together with the `--amcheck` flag when the `amcheck` extension is installed in the database.

The verification of unique constraints is only possible if in the version of the `amcheck` extension you are using, the `bt_index_check` function takes the `checkunique` parameter.

`--heapallindexed`

Checks that all heap tuples that should be indexed are actually indexed. You can use this flag only together with the `--amcheck` flag.

This check is only possible if in the version of the `amcheck/amcheck_next` extension you are using, the `bt_index_check` function takes the `heapallindexed` parameter.

`--skip-block-validation`

Skip validation of data files. You can use this flag only together with the `--amcheck` flag, so that only logical verification of indexes is performed.

For more details of the command settings, see sections [Common Options](#), [Connection Options](#), [Logging Options](#), and [S3 Options](#).

For details on usage, see the section [Verifying a Cluster](#).

### validate

```
pg_probackup validate -B backup_dir
[--help] [--instance=instance_name] [-i backup_id]
[-j num_threads] [--progress] [--wal]
[--skip-block-validation]
[recovery_target_options] [logging_options] [s3_options]
```

Verifies that all the files required to restore the cluster are present and are not corrupt. If `instance_name` is not specified, `pg_probackup` validates all backups available in the backup catalog. If you specify the `instance_name` without any additional options, `pg_probackup` validates all the backups available for this backup instance. If you specify the `instance_name` with a [recovery target option](#) and/or a `backup_id`, `pg_probackup` checks whether it is possible to restore the cluster using these options. If the `--wal` option is specified, full check of the WAL archive will be performed instead of only checking WAL segments needed to restore the cluster.

For details, see the section [Validating a Backup](#).

### Note

The `TATLIN.OBJECT` S3 storage is an “eventually consistent” class system. In such systems, the reconciliation across nodes proceeds asynchronously. Therefore, if you are making a backup to the `TATLIN.OBJECT` S3 storage, perform backup validation not immediately after the backup, but some time later.

#### merge

```
pg_probackup merge -B backup_dir --instance=instance_name -i backup_id
[--dry-run] [--help] [-j num_threads] [--progress] [--no-validate] [--no-sync]
[logging_options]
```

Merges backups that belong to a common incremental backup chain. If you specify a full backup, it will be merged with its first incremental backup. If you specify an incremental backup, it will be merged to its parent full backup, together with all incremental backups between them. Once the merge is complete, the full backup takes in all the merged data, and the incremental backups are removed as redundant.

`--no-validate`

Skips automatic validation before and after merge.

`--no-sync`

Do not sync merged files to disk. You can use this flag to speed up the merge process. Using this flag can result in data corruption in case of operating system or hardware crash.

For more details of the command settings, see sections [Common Options](#) and [Merging Backups](#).

#### delete

```
pg_probackup delete -B backup_dir --instance=instance_name
[--help] [-j num_threads] [--progress]
[--retention-redundancy=redundancy] [--retention-window=window] [--wal-depth=wal_depth]
[--delete-wal]
{-i backup_id | --delete-expired [--merge-expired] | --merge-expired | --
status=backup_status}
[--dry-run] [--no-validate] [--no-sync] [logging_options] [s3_options]
```

Deletes backup with specified *backup\_id* or launches the retention purge of backups and archived WAL that do not satisfy the current retention policies.

`--no-validate`

Skips automatic validation before and after retention merge.

`--no-sync`

Do not sync merged files to disk. You can use this flag to speed up the retention merge process. Using this flag can result in data corruption in case of operating system or hardware crash.

For details, see the sections [Deleting Backups](#), [Retention Options](#), and [Configuring Retention Policy](#).

#### archive-push

```
pg_probackup archive-push -B backup_dir --instance=instance_name
--wal-file-name=wal_file_name [--wal-file-path=wal_file_path]
[--help] [--dry-run] [--no-sync] [--compress] [--no-ready-rename] [--overwrite]
[-j num_threads] [--batch-size=batch_size]
[--archive-timeout=timeout]
[--compress-algorithm=compression_algorithm]
```

```
[--compress-level=compression_level]
[remote_options] [logging_options] [s3_options]
```

Copies WAL files into the corresponding subdirectory of the backup catalog and validates the backup instance by *instance\_name* and *system-identifier*. If parameters of the backup instance and the cluster do not match, this command fails with the following error message: Refuse to push WAL segment *segment\_name* into archive. Instance parameters mismatch.

If the files to be copied already exists in the backup catalog, `pg_probackup` computes and compares their checksums. If the checksums match, `archive-push` skips the corresponding file and returns a successful execution code. Otherwise, `archive-push` fails with an error. If you would like to replace WAL files in the case of checksum mismatch, run the `archive-push` command with the `--overwrite` flag.

Each file is copied to a temporary file with the `.part` suffix. If the temporary file already exists, `pg_probackup` will wait `archive_timeout` seconds before discarding it. After the copy is done, atomic rename is performed. This algorithm ensures that a failed `archive-push` will not stall continuous archiving and that concurrent archiving from multiple sources into a single WAL archive has no risk of archive corruption.

The Postgres Pro server requests WAL segments one at a time. To speed up archiving, you can specify the `--batch-size` option to copy WAL segments in batches of the specified size. If `--batch-size` option is used, then you can also specify the `-j` option to copy the batch of WAL segments on multiple threads.

WAL segments copied to the archive are synced to disk unless the `--no-sync` flag is used.

You can use `archive-push` in the [archive\\_command](#) Postgres Pro parameter to set up [continuous WAL archiving](#).

For more details of the command settings, see sections [Common Options](#), [Archiving Options](#), [Compression Options](#), and [S3 Options](#).

#### archive-get

```
pg_probackup archive-get -B backup_dir --instance=instance_name --wal-file-
path=wal_file_path --wal-file-name=wal_file_name
[-j num_threads] [--batch-size=batch_size]
[--prefetch-dir=prefetch_dir_path] [--no-validate-wal]
[--dry-run] [--help] [remote_options] [logging_options] [s3_options]
```

Copies WAL files from the corresponding subdirectory of the backup catalog to the cluster's write-ahead log location. This command is automatically set by `pg_probackup` as part of the `restore_command` when restoring backups using a WAL archive. You do not need to set it manually if you use local storage for backups or remote mode.

If you use S3 interface, to ensure that the Postgres Pro server has access to S3 storage to fetch WAL files during restore, you can specify the `--s3-config-file` option that defines the S3 configuration file with appropriate configuration settings, as described in [the section called “S3 Options”](#).

The Postgres Pro server requests WAL segments one at a time. To speed up recovery, you can specify the `--batch-size` option to copy WAL segments in batches of the specified size. If `--batch-size` option is used, then you can also specify the `-j` option to copy the batch of WAL segments on multiple threads.

For more details of the command settings, see sections [Common Options](#), [Archiving Options](#), [Compression Options](#), and [S3 Options](#).

#### catchup

```
pg_probackup catchup -b catchup_mode
--source-pgdata=path_to_pgdata_on_remote_server
--destination-pgdata=path_to_local_dir
[--help] [-j | --threads=num_threads] [--dry-run]
[--write-rate-limit=bitrate]
```

```
[--stream [--temp-slot[=true|false|on|off]] [-P | --perm-slot] [-S | --slot=slot_name]]
[--exclude-path=PATHNAME]
[-T OLDDIR=NEWDIR]
[-X | --waldir=wal_dir]
[connection_options] [remote_options]
[logging_options]
```

Creates a copy of a Postgres Pro instance without using the backup catalog.

`-b catchup_mode`

`--backup-mode=catchup_mode`

Specifies the catchup mode to use. Possible values are: [FULL](#), [DELTA](#), and [PTRACK](#).

`--source-pgdata=path_to_pgdata_on_remote_server`

Specifies the path to the data directory of the instance to be copied. The path can be local or remote.

`--destination-pgdata=path_to_local_dir`

Specifies the path to the local data directory to copy to.

`-j num_threads`

`--threads=num_threads`

Sets the number of parallel threads for `catchup` process.

`--stream`

Copies the instance in [STREAM](#) WAL delivery mode, including all the necessary WAL files by streaming them from the server via replication protocol. For `catchup`, it is enabled by default.

`--write-rate-limit=bitrate`

Sets the rate of writing data to disk, in MBps or GBps. The default unit is MBps. For example: `--write-rate-limit=1GBps` or `--write-rate-limit=100` (MBps). The default value is 0 — no limitation.

If this option is specified, the following information is displayed at the end of the `catchup` operation:

- `written` — the amount of data written, in MB.
- `total time` — the time that elapsed between the first and last writes, in seconds. Note that this is not the total duration of the `catchup` operation.
- `sleep time` — the amount of forced delay time, in seconds.
- `average rate` — the actual average write rate, in MBps.

For example:

```
INFO: Rate limit: written 14975.445 MB, total time 17.163 s, sleep time 2.370 s,
 average rate 872.560715 MBps
```

`-x=path_prefix`

`--exclude-path=path_prefix`

Specifies a prefix for files to exclude from the synchronization of Postgres Pro instances during copying. The prefix must contain a path relative to the data directory of an instance. If the prefix specifies a directory, all files in this directory will not be synchronized.

### Warning

This option is dangerous since excluding files from synchronization can result in incomplete synchronization; use with care.



`--temp-slot [=true|false|on|off]`

Creates a *temporary* physical replication slot for streaming WAL from the Postgres Pro instance being copied. `--temp-slot` is enabled by default. It ensures that all the required WAL segments remain available if WAL is rotated while the backup is in progress. This flag can only be used together with the `--stream` flag and cannot be used together with the `--perm-slot` flag. The default slot name is `pg_probackup_slot`. To change it, use the `--slot/-S` option and explicitly specify `--temp-slot` or `--temp-slot=true|on`.

`-P`

`--perm-slot`

Creates a *permanent* physical replication slot for streaming WAL from the Postgres Pro instance being copied. This flag can only be used together with the `--stream` flag and cannot be used together with the `--temp-slot` flag. The default slot name is `pg_probackup_perm_slot`, which can be changed using the `--slot/-S` option.

`-S slot_name`

`--slot=slot_name`

Specifies the replication slot to connect to for WAL streaming. This option can only be used together with the `--stream` flag.

`-T OLDDIR=NEWDIR`

`--tablespace-mapping=OLDDIR=NEWDIR`

Relocates the tablespace from the *OLDDIR* to the *NEWDIR* directory at the time of recovery. Both *OLDDIR* and *NEWDIR* must be absolute paths. If the path contains the equals sign (=), escape it with a backslash. This option can be specified multiple times for multiple tablespaces.

`-X wal_dir`

`--waldir=wal_dir`

Sets the directory to write WAL files to. By default WAL files will be placed in the `pg_wal` subdirectory of the target directory, but this option can be used to place them elsewhere. *wal\_dir* must be an absolute path, which must not already exist, but if it does, it must be empty to perform `catchup` with `--catchup-mode=FULL`.

For more details of the command settings, see sections [Common Options](#), [Connection Options](#), and [Remote Mode Options](#).

For details on usage, see the section [Cloning and Synchronizing Postgres Pro Instance](#).

## Options

This section describes command-line options for `pg_probackup` commands. If the option value can be derived from an environment variable, this variable is specified below the command-line option, in the uppercase. Some values can be taken from the `pg_probackup.conf` configuration file located in the backup catalog.

For details, see [the section called “Configuring pg\\_probackup”](#).

If an option is specified using more than one method, command-line input has the highest priority, while the `pg_probackup.conf` settings have the lowest priority.

## Common Options

The list of general options.

`--dry-run`

Initiates a trial run of the appropriate command, which does not actually do any changes, that is, it does not create, delete or move files on disk. This flag allows you to check that all the command options are correct and the command is ready to run. WAL streaming is skipped with `--dry-run`.



```
-B directory
--backup-path=directory
BACKUP_PATH
```

Specifies the absolute path to the backup catalog. Backup catalog is a directory where all backup files and meta information are stored. Since this option is required for most of the `pg_probackup` commands, you are recommended to specify it once in the `BACKUP_PATH` environment variable. In this case, you do not need to use this option each time on the command line.

```
-D directory
--pgdata=directory
PGDATA
```

Specifies the absolute path to the data directory of the database cluster. This option is mandatory only for the `add-instance` command. Other commands can take its value from the `PGDATA` environment variable, or from the `pg_probackup.conf` configuration file.

```
-i backup_id
--backup-id=backup_id
```

Specifies the unique identifier of the backup.

```
-j num_threads
--threads=num_threads
```

Sets the number of parallel threads for backup, restore, merge, validate, checkdb, and archive-push processes.

```
--progress
```

Shows the progress of operations.

```
--help
```

Shows detailed information about the options that can be used with this command.

### Recovery Target Options

If `continuous WAL archiving` is configured, you can use one of these options together with `restore` or `validate` commands to specify the moment up to which the database cluster must be restored or validated.

```
--recovery-target=immediate|latest
```

Defines when to stop the recovery:

- The `immediate` value stops the recovery after reaching the consistent state of the specified backup, or the latest available backup if the `-i/--backup-id` option is omitted. This is the default behavior for STREAM backups.
- The `latest` value continues the recovery until all WAL segments available in the archive are applied. Setting this value of `--recovery-target` also sets `--recovery-target-timeline` to `latest`.

```
--recovery-target-timeline=timeline
```

Specifies a particular timeline to be used for recovery:

- `current` — the timeline of the specified backup, default.
- `latest` — the timeline of the latest available backup.
- A numeric value.

```
--recovery-target-lsn=lsn
```

Specifies the LSN of the write-ahead log location up to which recovery will proceed.

`--recovery-target-name=recovery_target_name`

Specifies a named savepoint up to which to restore the cluster.

`--recovery-target-time=time`

Specifies the timestamp up to which recovery will proceed. If the time zone offset is not specified, the local time zone is used.

**Example:** `--recovery-target-time="2027-04-09 18:21:32+00"`

`--recovery-target-xid=xid`

Specifies the transaction ID up to which recovery will proceed.

`--recovery-target-inclusive=boolean`

Specifies whether to stop just after the specified recovery target (`true`), or just before the recovery target (`false`). This option can only be used together with `--recovery-target-time`, `--recovery-target-lsn` or `--recovery-target-xid` options. The default depends on the [recovery\\_target\\_inclusive](#) parameter.

`--recovery-target-action=pause|promote|shutdown`

Specifies [recovery\\_target\\_action](#) the server should take when the recovery target is reached.

Default: `pause`

### Retention Options

You can use these options together with [backup](#) and [delete](#) commands.

For details on configuring retention policy, see the section [Configuring Retention Policy](#).

`--retention-redundancy=redundancy`

Specifies the number of full backup copies to keep in the data directory. Must be a non-negative integer. The zero value disables this setting.

Default: `0`

`--retention-window=window`

Number of days of recoverability. Must be a non-negative integer. The zero value disables this setting.

Default: `0`

`--wal-depth=wal_depth`

Number of latest valid backups on every timeline that must retain the ability to perform PITR. Must be a non-negative integer. The zero value disables this setting.

Default: `0`

`--delete-wal`

Deletes WAL files that are no longer required to restore the cluster from any of the existing backups.

`--delete-expired`

Deletes backups that do not conform to the retention policy defined in the `pg_probackup.conf` configuration file.

`--merge-expired`

Merges the oldest incremental backup that satisfies the requirements of retention policy with its parent backups that have already expired.

`--dry-run`

Displays the current status of all the available backups, without deleting or merging expired backups, if any.

### Pinning Options

You can use these options together with [backup](#) and [set-backup](#) commands.

For details on backup pinning, see the section [Backup Pinning](#).

`--ttl=ttl`

Specifies the amount of time the backup should be pinned. Must be a non-negative integer. The zero value unpins the already pinned backup. Supported units: ms, s, min, h, d (s by default).

Example: `--ttl=30d`

`--expire-time=time`

Specifies the timestamp up to which the backup will stay pinned. Must be an ISO-8601 compliant timestamp. If the time zone offset is not specified, the local time zone is used.

Example: `--expire-time="2027-04-09 18:21:32+00"`

### Logging Options

You can use these options with any command.

`--no-color`

Disable coloring for console log messages of warning and error levels.

`--log-level-console=log_level`

Controls which message levels are sent to the console log. Valid values are `verbose`, `log`, `info`, `warning`, `error` and `off`. Each level includes all the levels that follow it. The later the level, the fewer messages are sent. The `off` level disables console logging.

Default: `info`

### Note

All console log messages are going to `stderr`, so the output of [show](#) and [show-config](#) commands does not mingle with log messages.

`--log-level-file=log_level`

Controls which message levels are sent to a log file. Valid values are `verbose`, `log`, `info`, `warning`, `error`, and `off`. Each level includes all the levels that follow it. The later the level, the fewer messages are sent. The `off` level disables file logging.

Default: `off`

`--log-filename=log_filename`

Defines the filenames of the created log files. The filenames are treated as a `strftime` pattern, so you can use %-escapes to specify time-varying filenames.

Default: `pg_probackup.log`

For example, if you specify the `pg_probackup-%u.log` pattern, `pg_probackup` generates a separate log file for each day of the week, with `%u` replaced by the corresponding decimal number: `pg_probackup-1.log` for Monday, `pg_probackup-2.log` for Tuesday, and so on.

This option takes effect if file logging is enabled by the `--log-level-file` option.

`--error-log-filename=error_log_filename`

Defines the filenames of log files for error messages only. The filenames are treated as a `strftime` pattern, so you can use %-escapes to specify time-varying filenames.

Default: none

For example, if you specify the `error-pg_probackup-%u.log` pattern, `pg_probackup` generates a separate log file for each day of the week, with `%u` replaced by the corresponding decimal number: `error-pg_probackup-1.log` for Monday, `error-pg_probackup-2.log` for Tuesday, and so on.

This option is useful for troubleshooting and monitoring.

`--log-directory=log_directory`

Defines the directory in which log files will be created. You must specify the absolute path. This directory is created lazily, when the first log message is written.

Note that the directory for log files is always created locally even if backups are created in the S3 storage. So be sure to pass a local path in `log_directory` when needed.

Default: `$BACKUP_PATH/log/`

`--log-format-console=log_format`

Defines the format of the console log. Only set from the command line. Note that you cannot specify this option in the `pg_probackup.conf` configuration file through the `set-config` command and that the `backup` command also treats this option specified in the configuration file as an error. Possible values are:

- `plain` — sets the plain-text format of the console log.
- `json` — sets the JSON format of the console log.

Default: `plain`

`--log-format-file=log_format`

Defines the format of log files used. Possible values are:

- `plain` — sets the plain-text format of log files.
- `json` — sets the JSON format of log files.

Default: `plain`

`--log-rotation-size=log_rotation_size`

Maximum size of an individual log file. If this value is reached, the log file is rotated once a `pg_probackup` command is launched, except `help` and `version` commands. The zero value disables size-based rotation. Supported units: `kB`, `MB`, `GB`, `TB` (`kB` by default).

Default: `0`

`--log-rotation-age=log_rotation_age`

Maximum lifetime of an individual log file. If this value is reached, the log file is rotated once a `pg_probackup` command is launched, except `help` and `version` commands. The time of the last log file creation is stored in `$BACKUP_PATH/log/log_rotation`. The zero value disables time-based rotation. Supported units: `ms`, `s`, `min`, `h`, `d` (`min` by default).

Default: `0`

**Connection Options**

You can use these options together with [backup](#), [catchup](#), and [checkdb](#) commands.

All [libpq environment variables](#) are supported.

```
-d dbname
--pgdatabase=dbname
PGDATABASE
```

Specifies the name of the database to connect to. The connection is used only for managing backup process, so you can connect to any existing database. If this option is not provided on the command line, PGDATABASE environment variable, or the `pg_probackup.conf` configuration file, `pg_probackup` tries to take this value from the PGUSER environment variable, or from the current user name if PGUSER variable is not set.

```
-h host
--pghost=host
PGHOST
```

Specifies the host name of the system on which the server is running. If the value begins with a slash, it is used as a directory for the Unix domain socket.

Default: localhost

```
-p port
--pgport=port
PGPORT
```

Specifies the TCP port or the local Unix domain socket file extension on which the server is listening for connections.

Default: 5432

```
-U username
--pguser=username
PGUSER
```

User name to connect as.

```
-w
--no-password
```

Disables a password prompt. If the server requires password authentication and a password is not available by other means such as a [pgpass](#) file or PGPASSWORD environment variable, the connection attempt will fail. This flag can be useful in batch jobs and scripts where no user is present to enter a password.

```
-W
--password
```

Forces a password prompt. (Deprecated)

**Compression Options**

You can use these options together with [backup](#) and [archive-push](#) commands.

```
--compress-algorithm=compression_algorithm
```

Defines the algorithm to use for compressing data files. Possible values are `zlib`, `lz4`, `zstd`, `pglz`, and `none`. If set to any value, but `none`, this option enables compression that uses the corresponding algorithm. Both data files and WAL files are compressed. By default, compression is disabled. For the [archive-push](#) command, the `pglz` compression algorithm is not supported.

**Note**

pg\_probackup supports compression algorithms included in the Postgres Pro version. In particular:

- lz4 is supported for Postgres Pro Enterprise 13 and higher.
- zstd is supported for Postgres Pro Enterprise 11 and higher.

Default: none

`--compress-level=compression_level`

Defines the compression level. This option can be used together with the `--compress-algorithm` option. Possible values depend on the compression algorithm specified:

- 0 — 9 for `zlib`
- 1 for `pglz`
- 0 — 12 for `lz4`
- 0 — 22 for `zstd`

The value of 0 sets the default compression level for the specified algorithm:

- 6 for `zlib`
- 1 for `pglz`
- 9 for `lz4`
- 3 for `zstd`

**Note**

The pure lz4 algorithm has only one compression level — 1. So, if the specified compression algorithm is `lz4` and `--compress-level` is greater than 1, the `lz4hc` algorithm is actually used, which is much slower although does better compression.

Default: 1

`--compress`

Specifies the default compression algorithm and `--compress-level=1`. The default algorithm is selected among those supported by Postgres Pro according to the priorities: `zstd` (highest) -> `lz4` -> `zlib` -> `pglz`. The `--compress` option overrides the `--compression-algorithm` and `--compress-level` settings and cannot be specified together with them.

**Archiving Options**

These options can be used with the [archive-push](#) command in the [archive\\_command](#) setting and the [archive-get](#) command in the [restore\\_command](#) setting.

Additionally, [remote mode options](#) and [logging options](#) can be used.

`--wal-file-path=wal_file_path`

Provides the path to the WAL file in `archive_command` and `restore_command`. Use the `%p` variable as the value for this option or explicitly specify the path to a file outside of the data directory. If you skip this option, the path specified in `pg_probackup.conf` will be used.

`--wal-file-name=wal_file_name`

Provides the name of the WAL file in `archive_command` and `restore_command`. Use the `%f` variable as the value for this option for correct processing. If the value of `--wal-file-path` is a path outside of the data directory, explicitly specify the filename.

`--overwrite`

Overwrites archived WAL file. Use this flag together with the [archive-push](#) command if the specified subdirectory of the backup catalog already contains this WAL file and it needs to be replaced with its newer copy. Otherwise, `archive-push` reports that a WAL segment already exists, and aborts the operation. If the file to replace has not changed, `archive-push` skips this file regardless of the `--overwrite` flag.

`--batch-size=batch_size`

Used to speed up archiving in case of `archive-push` or to speed up recovery in case of `archive-get`. Sets the maximum number of WAL files that can be copied into the archive by a single `archive-push` process, or from the archive by a single `archive-get` process.

`--archive-timeout=wait_time`

Sets the timeout for considering existing `.part` files to be stale. By default, `pg_probackup` waits 300 seconds. This option can be used only with [archive-push](#) command.

`--no-ready-rename`

Do not rename status files in the `archive_status` directory. This option should be used only if `archive_command` contains multiple commands. This option can be used only with [archive-push](#) command.

`--no-sync`

Do not sync copied WAL files to disk. You can use this flag to speed up archiving process. Using this flag can result in WAL archive corruption in case of operating system or hardware crash. This option can be used only with [archive-push](#) command.

`--prefetch-dir=path`

Directory used to store prefetched WAL segments if `--batch-size` option is used. Directory must be located on the same filesystem and on the same mountpoint the `PGDATA/pg_wal` is located. By default files are stored in `PGDATA/pg_wal/pbk_prefetch` directory. This option can be used only with [archive-get](#) command.

`--no-validate-wal`

Do not validate prefetched WAL file before using it. Use this option if you want to increase the speed of recovery. This option can be used only with [archive-get](#) command.

### Remote Mode Options

This section describes the options related to running `pg_probackup` operations remotely via SSH. These options can be used with [add-instance](#), [set-config](#), [backup](#), [catchup](#), [restore](#), [archive-push](#), and [archive-get](#) commands.

For details on configuring and using the remote mode, see [the section called “Configuring the Remote Mode”](#) and [the section called “Using pg\\_probackup in the Remote Mode”](#).

`--remote-proto=proto`

Specifies the protocol to use for remote operations. Currently only the SSH protocol is supported. Possible values are:

- `ssh` enables the remote mode via SSH. This is the default value.

- `none` explicitly disables the remote mode.

You can omit this option if the `--remote-host` option is specified.

`--remote-host=destination`

Specifies the remote host IP address or hostname to connect to.

`--remote-port=port`

Specifies the remote host port to connect to.

Default: 22

`--remote-user=username`

Specifies remote host user for SSH connection. If you omit this option, the current user initiating the SSH connection is used.

`--remote-path=path`

Specifies `pg_probackup` installation directory on the remote system.

`--ssh-options=ssh_options`

Provides a string of SSH command-line options. For example, the following options can be used to set `keep-alive` for SSH connections opened by `pg_probackup`: `--ssh-options="-o ServerAliveCountMax=5 -o ServerAliveInterval=60"`. For the full list of possible options, see [ssh\\_config manual page](#).

### Remote WAL Archive Options

This section describes the options used to provide the arguments for [remote mode options](#) in `archive-get` used in the `restore_command` command when restoring ARCHIVE backups or performing PITR.

`--archive-host=destination`

Provides the argument for the `--remote-host` option in the `archive-get` command.

`--archive-port=port`

Provides the argument for the `--remote-port` option in the `archive-get` command.

Default: 22

`--archive-user=username`

Provides the argument for the `--remote-user` option in the `archive-get` command. If you omit this option, the user that has started the Postgres Pro cluster is used.

Default: Postgres Pro user

### Incremental Restore Options

This section describes the options for incremental cluster restore. These options can be used with the `restore` command.

`-I incremental_mode`

`--incremental-mode=incremental_mode`

Specifies the incremental mode to be used. Possible values are:

- `CHECKSUM` — replace only pages with mismatched checksum and LSN.
- `LSN` — replace only pages with LSN greater than point of divergence.
- `NONE` — regular restore.



### Partial Restore Options

This section describes the options for partial cluster restore. These options can be used with the [restore](#) command.

`--db-exclude=dbname`

Specifies the name of the database to exclude from restore. All other databases in the cluster will be restored as usual, including `template0` and `template1`. This option can be specified multiple times for multiple databases.

`--db-include=dbname`

Specifies the name of the database to restore from a backup. All other databases in the cluster will not be restored, with the exception of `template0` and `template1`. This option can be specified multiple times for multiple databases.

### S3 Options

This section describes the options needed to store backups in private clouds. These options can be used with any commands that `pg_probackup` runs using S3 interface.

`--s3=s3_interface_provider`

Specifies the S3 interface provider. Possible values are:

- `minio` — MinIO object storage, compatible with S3 cloud storage service. With this provider, custom S3 server settings can be specified. The HTTP protocol, port 9000, and region `us-east-1` are used by default.
- `vk` — VK Cloud storage. With this provider, the S3 host address `hb.vkcs.cloud`, port 443, and HTTPS protocol are only used. Custom values of the host, port, and protocol are ignored. The default value of region is `ru-msk`.
- `aws` — Amazon S3 storage, offered by Amazon Web Services (AWS). With this provider, the S3 host address `bucket_name.s3.region.amazonaws.com`, port 443, and HTTPS protocol are only used. Custom values of the host, port, and protocol are ignored. The default value of region is `us-east-1`.

With `--s3=minio`, `pg_probackup` will work fine for a VK Cloud storage if the S3 host address, port and protocol are properly specified (host address is `hb.vkcs.cloud` or the one specified in the appropriate section of the VK Cloud profile, port 443, and HTTPS protocol). Do not specify `--s3=minio` for the Amazon S3 storage.

Once a `pg_probackup` command runs with the `--s3` option, `pg_probackup` starts running all commands that support parallel execution on 10 parallel threads (for details, see [the section called “Running pg\\_probackup on Parallel Threads”](#)). You can change the number of threads using the `-j/--threads` option.

`--s3-config-file=path_to_config_file`

Specifies the S3 configuration file. Settings in the configuration file override the environment variables. If this option is not specified, `pg_probackup` first looks for the S3 configuration file at `/etc/pg_probackup/s3.config` and then at `~postgres/.pg_probackup/s3.config`. The following is an example of the S3 configuration file:

```
access-key = ...
secret-key = ...
s3-host = localhost
s3-port = 9000
s3-bucket = s3demo
s3-region=us-east-1
s3-buffer-size = 32
s3-secure = on | https | http | off
```

## Testing and Debugging Options

This section describes options useful only in a test or development environment.

`--cfs-nondatfile-mode`

Instructs `backup` command to backup CFS in a legacy mode. This allows fine-tuning compatibility with `pg_probackup` versions earlier than 2.6.0. This option is mainly designed for testing.

`PGPROBACKUP_TESTS_SKIP_HIDDEN`

Instructs `pg_probackup` to ignore backups marked as hidden. Note that `pg_probackup` can never mark a backup as hidden. It can only be done by directly editing the `backup.control` file. This option can only be set with environment variables.

`--destroy-all-other-dbs`

By default, `pg_probackup` exits with an error if an attempt is made to perform a partial incremental restore since this destroys databases not included in the restore set. This flag allows you to suppress the error and proceed with the partial incremental restore (e.g., to keep a development database snapshot up-to-date with a production one). This option can be used with the `restore` command.

### Important

Never use this flag in a production cluster.

`--lock-lifetime`

Specifies the lifetime for locks taken when performing a backup, in seconds. The default value is 300. The minimum value is 60.

`PGPROBACKUP_TESTS_SKIP_EMPTY_COMMIT`

Instructs `pg_probackup` to skip empty commits after `pg_backup_stop`.

## Versioning

`pg_probackup` follows *semantic* versioning.

## Authors

Postgres Professional, Moscow, Russia.

## Credits

`pg_probackup` utility is based on `pg_arman`, which was originally written by NTT and then developed and maintained by Michael Paquier.

## vacuumlo

`vacuumlo` — remove orphaned large objects from a Postgres Pro database

## Synopsis

```
vacuumlo [option...] dbname...
```

## Description

`vacuumlo` is a simple utility program that will remove any “orphaned” large objects from a Postgres Pro database. An orphaned large object (LO) is considered to be any LO whose OID does not appear in any `oid` or `lo` data column of the database.

All databases named on the command line are processed.

## Options

`vacuumlo` accepts the following command-line arguments:

```
-l limit
--limit=limit
```

Remove no more than *limit* large objects per transaction (default 1000). Since the server acquires a lock per LO removed, removing too many LOs in one transaction risks exceeding [max\\_locks\\_per\\_transaction](#). Set the limit to zero if you want all removals done in a single transaction.

```
-n
--dry-run
```

Don't remove anything, just show what would be done.

```
-v
--verbose
```

Write a lot of progress messages.

```
-V
--version
```

Print the `vacuumlo` version and exit.

```
-?
--help
```

Show help about `vacuumlo` command line arguments, and exit.

`vacuumlo` also accepts the following command-line arguments for connection parameters:

```
-h host
--host=host
```

Database server's host.

```
-p port
--port=port
```

Database server's port.

```
-U username
--username=username
```

User name to connect as.

```
-w
--no-password
```

Never issue a password prompt. If the server requires password authentication and a password is not available by other means such as a `.pgpass` file, the connection attempt will fail. This option can be useful in batch jobs and scripts where no user is present to enter a password.

```
-W
--password
```

Force `vacuumlo` to prompt for a password before connecting to a database.

This option is never essential, since `vacuumlo` will automatically prompt for a password if the server demands password authentication. However, `vacuumlo` will waste a connection attempt finding out that the server wants a password. In some cases it is worth typing `-w` to avoid the extra connection attempt.

## Environment

```
PGHOST
PGPORT
PGUSER
```

Default connection parameters.

This utility, like most other Postgres Pro utilities, also uses the environment variables supported by `libpq` (see [Section 35.15](#)).

The environment variable `PG_COLOR` specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

## Notes

`vacuumlo` works by the following method: First, `vacuumlo` builds a temporary table which contains all of the OIDs of the large objects in the selected database. It then scans through all columns in the database that are of type `oid` or `lo`, and removes matching entries from the temporary table. (Note: Only types with these names are considered; in particular, domains over them are not considered.) The remaining entries in the temporary table identify orphaned LOs. These are removed.

## Author

Peter Mount <[peter@retep.org.uk](mailto:peter@retep.org.uk)>

## H.2. Additional Postgres Pro Server Applications

This section covers additional Postgres Pro Shardman server-related applications. They are typically run on the host where the database server resides. See also [Postgres Pro Server Applications](#) for information about server applications that are part of the core Postgres Pro Shardman distribution.

## pgpro\_controldata

`pgpro_controldata` — display control information of a PostgreSQL/Postgres Pro database cluster and compatibility information for a cluster and/or server

### Synopsis

```
pgpro_controldata [option...]
```

### Description

`pgpro_controldata` prints control information, such as the catalog version, initialized by `initdb` command of any PostgreSQL/Postgres Pro server. It also shows information about write-ahead logging and checkpoint processing. This information is cluster-wide and not specific to any database.

`pgpro_controldata` also helps to check compatibility between PostgreSQL/Postgres Pro database servers and clusters. It can print server or cluster parameters that can affect the compatibility and check whether a cluster and server are compatible by comparing those parameters.

`pgpro_controldata` is provided with Postgres Pro Shardman as a separate pre-built package `pgpro-controldata` (for the detailed installation instructions, see [Chapter 17](#)).

### Options

`pgpro_controldata` accepts the following command-line arguments. If no arguments are specified, `pgpro_controldata` just prints the control information. Note that compatibility-related command-line arguments `-P` and `-S` specified together work the same way as `-C`.

#### General-Purpose

`-B`  
`--bindir`  
Specifies the PostgreSQL/Postgres Pro executable directory, needed to get server compatibility parameters.

`-D datadir`  
`--pgdata=datadir`  
Specifies the file system location of the database configuration files. If this option is omitted, the environment variable `PGDATA` is used.

`-V`  
`--version`  
Print the `pgpro_controldata` version, then exit.

`-?`  
`--help`  
Show help about `pgpro_controldata` command-line arguments, then exit.

#### Compatibility-Related

`-C`  
`--compatibility-check`  
Display all parameters that can affect compatibility between the specified server and cluster and check whether they are compatible.

Use the `-D` option or the environment variable `PGDATA` to provide the path to the data directory, where read access is required.

If the `-B` option is omitted, the current server is assumed.

The cluster data and the server must have the same byte order and architecture type for this option to work correctly.

`-P`

`--cluster-compatibility-params`

Display all parameters of the specified cluster that can affect the compatibility.

Use the `-D` option or the environment variable `PGDATA` to provide the path to the data directory, where read access is required.

The cluster data and the server must have the same byte order and architecture type for this option to work correctly.

`-S`

`--server-compatibility-params`

Display all parameters of the specified or current server that can affect the compatibility.

If the `-B` option is omitted, the current server is assumed.

## Environment

`PGDATA`

Default data directory location

`PG_COLOR`

Specifies whether to use color in diagnostic messages. Possible values are `always`, `auto` and `never`.

## pg\_integrity\_check

`pg_integrity_check` — calculate and validate checksums for controlled files (certified edition only)

### Synopsis

```
pg_integrity_check [connection-option...] [-s | --system] [-u | --user] [-c | --catalog] [-o | --output]
[-l filename | --log=filename] [--syslog] [-D datadir] [-C filename] [-v | --verbose] [-? | --help]
```

### Description

`pg_integrity_check` is a utility provided with Postgres Pro Shardman that can calculate and validate checksums for the objects you would like to control. When running `pg_integrity_check`, you must specify at least one of the options that define the type of the controlled objects: `-s`, `-u`, or `-c` for read-only files, additional files, or system catalog tables, respectively.

If you use the `-o` option, `pg_integrity_check` calculates checksums and writes them into configuration files under the `share/security/` directory. You must have write access to the corresponding files to perform this command. You cannot use the `-o` option together with the `-s` option, since a checksum for read-only files cannot be overwritten.

If you omit the `-o` option, `pg_integrity_check` compares the calculated checksums with the corresponding checksums in the configuration files. If the checksums differ for any of the controlled objects, `pg_integrity_check` displays a message indicating the difference.

For details on using `pg_integrity_check`, see [Section 31.2](#).

### Options

*connection-options*

Standard options for connecting to a database: `-d`, `-h`, `-p`, `-U`. You must specify `-d` and `-U` options when validating checksums for catalog tables with the `-c` option. For detailed description of connection options, see [psql](#).

`-s`  
`--system`

Validate checksums for read-only files. Checksums for read-only files control both file contents and file attributes.

`-u`  
`--user`

Validate checksums for additional files. Checksums for additional files control both file contents and file attributes.

`-c`  
`--catalog`

Validate checksums for system catalog tables. For the `-c` option to work correctly, you must also specify connection parameters for the database. The database server must be started and accept connections.

`-o`  
`--output`

Recalculate checksums and write it into a file.

`-l filename`  
`--log=filename`

Write checksum validation results into a log file.

`--syslog`

Write checksum validation results into the syslog.

`-D datadir`

Data directory of the database cluster. This option is required to define the filenames for the generated configuration files when using `-u` option.

`-C filename`

The absolute path to the configuration file used for integrity checks of system catalog tables of the selected database. If used with the `-o` option, `pg_integrity_check` writes into the specified file. If no other option is used, `pg_integrity_check` checks data against the specified file.

`-v`

`--verbose`

Print debugging information, including checksum values.

`-?`

`--help`

Print command-line help.

## Return Values

0 — checksums are calculated or validated successfully.

1 — an error occurred during the initial checksum calculation for read-only files.

2 — checksum validation revealed changes in one or more of the controlled objects.

3 — an unexpected error occurred during checksum validation.

## Examples

Compute checksums for additional files and write them into the `share/security/_var_lib_pgpro_sdm-17_data.user.conf` configuration file:

```
pg_integrity_check -u -o -D /var/lib/pgpro/sdm-17/data
```

Check integrity of all controlled objects in the `postgres` database on behalf of the `postgres` user:

```
pg_integrity_check -s -u -c -D /var/lib/pgpro/sdm-17/data -d postgres -h localhost -p 5432 -U postgres
```

## H.3. Third-Party Server Applications

This section covers third-party server-related applications. They are typically run on the host where the database server resides.



## pgbouncer

pgbouncer — a Postgres Pro connection pooler

### Synopsis

On Linux systems:

```
pgbouncer [-d] [-R] [-v] [-u user] pgbouncer.ini
```

```
pgbouncer -V | -h
```

On Windows:

```
pgbouncer [-v] [-u user] pgbouncer.ini
```

```
pgbouncer -V | -h
```

To use pgbouncer as a Windows service:

```
pgbouncer.exe --regservice pgbouncer.ini
```

```
pgbouncer.exe --unregservice pgbouncer.ini
```

### Description

pgbouncer is a Postgres Pro connection pooler. Any target application can be connected to pgbouncer as if it were a Postgres Pro server, and pgbouncer will create a connection to the actual server, or it will reuse one of its existing connections.

The aim of pgbouncer is to lower the performance impact of opening new connections to Postgres Pro.

In order not to compromise transaction semantics for connection pooling, pgbouncer supports several types of pooling when rotating connections:

#### Session pooling

Most polite method. When a client connects, a server connection will be assigned to it for the whole duration the client stays connected. When the client disconnects, the server connection will be put back into the pool. This is the default method.

#### Transaction pooling

A server connection is assigned to a client only during a transaction. When pgbouncer notices that transaction is over, the server connection will be put back into the pool.

#### Statement pooling

Most aggressive method. The server connection will be put back into the pool immediately after a query completes. Multi-statement transactions are disallowed in this mode as they would break.

The administration interface of pgbouncer consists of some new `SHOW` commands available when connected to a special “virtual” database `pgbouncer`.

### Quick Start

pgbouncer is provided with Postgres Pro Shardman as a separate pre-built package `pgbouncer` (for the detailed installation instructions, see [Chapter 17](#)). Basic setup and usage is as follows.

1. Create a `pgbouncer.ini` file. Details in the `pgbouncer(5)` man page. Simple example:

```
[databases]
template1 = host=localhost dbname=template1 auth_user=someuser
```

```
[pgbouncer]
listen_port = 6432
listen_addr = localhost
auth_type = md5
auth_file = userlist.txt
logfile = pgbouncer.log
pidfile = pgbouncer.pid
admin_users = someuser
```

2. Create a `userlist.txt` file that contains the users allowed in:

```
"someuser" "same_password_as_in_server"
```

3. Launch `pgbouncer`:

```
$ pgbouncer -d pgbouncer.ini
```

### Note

The above command does not work on Windows systems. Instead, `pgbouncer` must be launched as a service that first needs to be registered, as follows:

```
pgbouncer --regservice
```

4. Have your application (or the `psql` client) connect to `pgbouncer` instead of directly to the PostgreSQL server:

```
$ psql -p 6432 -U someuser template1
```

5. Manage `pgbouncer` by connecting to the special administration database `pgbouncer` and issuing `SHOW HELP`; to begin:

```
$ psql -p 6432 -U someuser pgbouncer
pgbouncer=# SHOW HELP;
NOTICE: Console usage
DETAIL:
SHOW [HELP|CONFIG|DATABASES|FDS|POOLS|CLIENTS|SERVERS|SOCKETS|LISTS|VERSION|...]
SET key = arg
RELOAD
PAUSE
SUSPEND
RESUME
SHUTDOWN
[...]
```

6. If you made changes to the `pgbouncer.ini` file, you can reload it with:

```
pgbouncer=# RELOAD;
```

## Options

`-d, --daemon`

Run in the background. Without it, the process will run in the foreground. In daemon mode, setting `pidfile` as well as `logfile` or `syslog` is required. No log messages will be written to `stderr` after going into the background.

### Note

Does not work on Windows, `pgbouncer` needs to run as service there.

`-R, --reboot`

### Note

This option is deprecated. Instead of this option use a rolling restart with multiple pgbouncer processes listening on the same port using `so_reuseport` instead.

Do an online restart. That means connecting to the running process, loading the open sockets from it, and then using them. If there is no active process, boot normally.

### Note

Works only if OS supports Unix sockets and the `unix_socket_dir` is not disabled in configuration. Does not work on Windows. Does not work with TLS connections, they are dropped.

`-u user, --user user`

Switch to the given user on startup.

`-v, --verbose`

Increase verbosity. Can be used multiple times.

`-q, --quiet`

Be quiet: do not log to stderr. This does not affect logging verbosity, only that stderr is not to be used. For use in `init.d` scripts.

`-V, --version`

Show version.

`-h, --help`

Show short help.

`--regservice`

Win32: Register to run as Windows service. The `service_name` configuration parameter value is used as the name to register under.

`--unregservice`

Win32: Unregister Windows service.

## Admin Console

The console is available by connecting as normal to the database pgbouncer:

```
$ psql -p 6432 pgbouncer
```

Only users listed in the configuration parameters `admin_users` or `stats_users` are allowed to log in to the console. (Except when `auth_mode=any`, then any user is allowed in as a `stats_user`.)

Additionally, the user name `pgbouncer` is allowed to log in without password, if the login comes via the Unix socket and the client has same Unix user uid as the running process.

The admin console currently only supports the simple query protocol. Some drivers use the extended query protocol for all commands; these drivers will not work for this.

## Show Commands

The `SHOW` commands output information. Each command is described below.

**SHOW STATS**

Shows statistics. In this and related commands, the total figures are since process start, the averages are updated every `stats_period`.

database

Statistics are presented per database.

total\_xact\_count

Total number of SQL transactions pooled by pgbouncer.

total\_query\_count

Total number of SQL queries pooled by pgbouncer.

total\_received

Total volume in bytes of network traffic received by pgbouncer.

total\_sent

Total volume in bytes of network traffic sent by pgbouncer.

total\_xact\_time

Total number of microseconds spent by pgbouncer when connected to Postgres Pro in a transaction, either idle in transaction or executing queries.

total\_query\_time

Total number of microseconds spent by pgbouncer when actively connected to Postgres Pro, executing queries.

total\_wait\_time

Time spent by clients waiting for a server, in microseconds. Updated when a client connection is assigned a backend connection.

avg\_xact\_count

Average transactions per second in last stat period.

avg\_query\_count

Average queries per second in last stat period.

avg\_recv

Average received (from clients) bytes per second.

avg\_sent

Average sent (to clients) bytes per second.

avg\_xact\_time

Average transaction duration, in microseconds.

avg\_query\_time

Average query duration, in microseconds.

avg\_wait\_time

Average time spent by clients waiting for a server that were assigned a backend connection within the current `stats_period`, in microseconds (average per second within that period).

### **SHOW STATS\_TOTALS**

Subset of `SHOW STATS` showing the total values (`total_`).

### **SHOW STATS\_AVERAGES**

Subset of `SHOW STATS` showing the average values (`avg_`).

### **SHOW TOTALS**

Like `SHOW STATS` but aggregated across all databases.

### **SHOW SERVERS**

`type`

S, for server.

`user`

User name pgbouncer uses to connect to server.

`database`

Database name.

`state`

State of the pgbouncer server connection, one of `active`, `idle`, `used`, `tested`, `new`, `active_cancel`, or `being_canceled`.

`addr`

IP address of Postgres Pro server.

`port`

Port of Postgres Pro server.

`local_addr`

Connection start address on local machine.

`local_port`

Connection start port on local machine.

`connect_time`

When the connection was made.

`request_time`

When last request was issued.

`wait`

Not used for server connections.

`wait_us`

Not used for server connections.

`close_needed`

1 if the connection will be closed as soon as possible, because a configuration file reload or DNS update changed the connection information or `RECONNECT` was issued.

`ptr`

Address of internal object for this connection. Used as unique ID.

link

Address of client connection the server is paired with.

remote\_pid

PID of backend server process. In case connection is made over Unix socket and OS supports getting process ID info, its OS PID. Otherwise it's extracted from cancel packet the server sent, which should be the PID in case the server is Postgres Pro, but it's a random number in case the server is another pgbouncer.

tls

A string with TLS connection information, or empty if not using TLS.

application\_name

A string containing the `application_name` set on the linked client connection, or empty if this is not set, or if there is no linked connection.

#### **SHOW CLIENTS**

type

C, for client.

user

Client connected user.

database

Database name.

state

State of the client connection, one of `active`, `waiting`, `active_cancel_req`, or `waiting_cancel_req`.

addr

IP address of the client.

port

Source port of the client.

local\_addr

Connection end address on local machine.

local\_port

Connection end port on local machine.

connect\_time

Timestamp of connect time.

request\_time

Timestamp of latest client request.

wait

Current waiting time in seconds.

wait\_us

Microsecond part of the current waiting time.

`close_needed`

Not used for clients.

`ptr`

Address of internal object for this connection. Used as unique ID.

`link`

Address of server connection the client is paired with.

`remote_pid`

Process ID, in case client connects over Unix socket and OS supports getting it.

`tls`

A string with TLS connection information, or empty if not using TLS.

`application_name`

A string containing the `application_name` set by the client for this connection, or empty if this is not set.

#### **SHOW POOLS**

A new pool entry is made for each couple of (database, user).

`database`

Database name.

`user`

User name.

`cl_active`

Client connections that are either linked to server connections or are idle with no queries waiting to be processed.

`cl_waiting`

Client connections that have sent queries but have not yet got a server connection.

`cl_active_cancel_req`

Client connections that have forwarded query cancellations to the server and are waiting for the server response.

`cl_waiting_cancel_req`

Client connections that have not forwarded query cancellations to the server yet.

`sv_active`

Server connections that are linked to a client.

`sv_active_cancel`

Server connections that are currently forwarding a cancel request.

`sv_being_canceled`

Servers that normally could become idle but are waiting to do so until all in-flight cancel requests have completed that were sent to cancel a query on this server.

sv\_idle

Server connections that are unused and immediately usable for client queries.

sv\_used

Server connections that have been idle for more than `server_check_delay`, so they need `server_check_query` to run on them before they can be used again.

sv\_tested

Server connections that are currently running either `server_reset_query` or `server_check_query`.

sv\_login

Server connections currently in the process of logging in.

maxwait

How long the first (oldest) client in the queue has waited, in seconds. If this starts increasing, then the current pool of servers does not handle requests quickly enough. The reason may be either an overloaded server or just too small of a `pool_size` setting.

maxwait\_us

Microsecond part of the maximum waiting time.

pool\_mode

The pooling mode in use.

#### **SHOW PEER\_POOLS**

A new `peer_pool` entry is made for each configured peer.

database

ID of the configured peer entry.

cl\_active\_cancel\_req

Client connections that have forwarded query cancellations to the server and are waiting for the server response.

cl\_waiting\_cancel\_req

Client connections that have not forwarded query cancellations to the server yet.

sv\_active\_cancel

Server connections that are currently forwarding a cancel request.

sv\_login

Server connections currently in the process of logging in.

#### **SHOW LISTS**

Show following internal information, in columns (not rows):

databases

Count of databases.

users

Count of users.



`pools`

Count of pools.

`free_clients`

Count of free clients.

`used_clients`

Count of used clients.

`login_clients`

Count of clients in `login` state.

`free_servers`

Count of free servers.

`used_servers`

Count of used servers.

`dns_names`

Count of DNS names in the cache.

`dns_zones`

Count of DNS zones in the cache.

`dns_queries`

Count of in-flight DNS queries.

`dns_pending`

Not used.

#### **SHOW USERS**

`name`

The user name.

`pool_mode`

The user's override `pool_mode`, or `NULL` if the default will be used instead.

#### **SHOW DATABASES**

`name`

Name of configured database entry.

`host`

Host pgbouncer connects to.

`port`

Port pgbouncer connects to.

`database`

Actual database name pgbouncer connects to.

`force_user`

When the user is part of the connection string, the connection between pgbouncer and Postgres Pro is forced to the given user, whatever the client user.

`pool_size`

Maximum number of server connections.

`min_pool_size`

Minimum number of server connections.

`reserve_pool`

Maximum number of additional connections for this database.

`pool_mode`

The database's override `pool_mode`, or `NULL` if the default will be used instead.

`max_connections`

Maximum number of allowed connections for this database, as set by `max_db_connections`, either globally or per database.

`current_connections`

Current number of connections for this database.

`paused`

1 if this database is currently paused, else 0.

`disabled`

1 if this database is currently disabled, else 0.

#### **SHOW PEERS**

`peer_id`

ID of the configured peer entry.

`host`

Host pgbouncer connects to.

`port`

Port pgbouncer connects to.

`pool_size`

Maximum number of server connections that can be made to this peer.

#### **SHOW FDS**

Internal command — shows list of file descriptors (FDs) in use with internal state attached to them.

When the connected user has the user name `pgbouncer`, connects through the Unix socket and has the same UID as the running process, the actual FDs are passed over the connection. This mechanism is used to do an online restart.

### **Note**

This does not work on Windows.

This command also blocks the internal event loop, so it should not be used while pgbouncer is in use.

`fd`

File descriptor numeric value.

`task`

One of `pooler`, `client` or `server`.

`user`

User of the connection using the FD.

`database`

Database of the connection using the FD.

`addr`

IP address of the connection using the FD, `unix` if a Unix socket is used.

`port`

Port used by the connection using the FD.

`cancel`

Cancel key for this connection.

`link`

File descriptor for corresponding server/client. `NULL` if idle.

#### **SHOW SOCKETS, SHOW ACTIVE\_SOCKETS**

Shows low-level information about sockets or only active sockets. This includes the information shown under `SHOW CLIENTS` and `SHOW SERVERS` as well as other more low-level information.

#### **SHOW CONFIG**

Show the current configuration settings, one per row, with the following columns:

`key`

Configuration variable name.

`value`

Configuration value.

`default`

Configuration default value.

`changeable`

Either `yes` or `no`, shows if the variable can be changed while running. If `no`, the variable can be changed only at boot-time. Use `SET` to change a variable at run time.

#### **SHOW MEM**

Shows low-level information about the current sizes of various internal memory allocations. The information presented is subject to change.

#### **SHOW DNS\_HOSTS**

Show host names in DNS cache.

hostname

Host name.

ttl

How many seconds until next lookup.

addrs

Comma separated list of addresses.

#### **SHOW DNS\_ZONES**

Show DNS zones in cache.

zonename

Zone name.

serial

Current serial.

count

Host names belonging to this zone.

#### **SHOW VERSION**

Show the pgbouncer version string.

#### **SHOW STATE**

Show the pgbouncer state settings. Current states are active, paused and suspended.

### **Process Controlling Commands**

#### **PAUSE [db]**

pgbouncer tries to disconnect from all servers. Disconnecting each server connection waits for that server connection to be released according to the server pool's pooling mode (in transaction pooling mode, the transaction must complete, in statement mode, the statement must complete, and in session pooling mode the client must disconnect). The command will not return before all server connections have been disconnected. To be used at the time of database restart.

If database name is given, only that database will be paused.

New client connections to a paused database will wait until `RESUME` is called.

#### **DISABLE db**

Reject all new client connections on the given database.

#### **ENABLE db**

Allow new client connections after a previous `DISABLE` command.

#### **RECONNECT db**

Close each open server connection for the given database, or all databases, after it is released (according to the pooling mode), even if its lifetime is not up yet. New server connections can be made immediately and will connect as necessary according to the pool size settings.

This command is useful when the server connection setup has changed, for example to perform a gradual switchover to a new server. It is not necessary to run this command when the connection string in `pgbouncer.ini` has been changed and reloaded (see `RELOAD`) or when DNS resolution has changed,

because then the equivalent of this command will be run automatically. This command is only necessary if something downstream of pgbouncer routes the connections.

After this command is run, there could be an extended period where some server connections go to an old destination and some server connections go to a new destination. This is likely only sensible when switching read-only traffic between read-only replicas, or when switching between nodes of a multimas-ter replication setup. If all connections need to be switched at the same time, `PAUSE` is recommended instead. To close server connections without waiting (for example, in emergency failover rather than gradual switchover scenarios), also consider `KILL`.

#### **KILL *db***

Immediately drop all client and server connections on given database.

New client connections to a killed database will wait until `RESUME` is called.

#### **SUSPEND**

All socket buffers are flushed and pgbouncer stops listening for data on them. The command will not return before all buffers are empty. To be used at the time of pgbouncer online reboot.

New client connections to a suspended database will wait until `RESUME` is called.

#### **RESUME [*db*]**

Resume work from previous `KILL`, `PAUSE`, or `SUSPEND` command.

#### **SHUTDOWN**

The pgbouncer process will exit.

#### **RELOAD**

The pgbouncer process will reload its configuration files and update changeable settings. This includes the main configuration file as well as the files specified by the settings `auth_file` and `auth_hba_file`.

pgbouncer notices when a configuration file reload changes the connection parameters of a database definition. An existing server connection to the old destination will be closed when the server connection is next released (according to the pooling mode), and new server connections will immediately use the updated connection parameters.

#### **WAIT\_CLOSE [*db*]**

Wait until all server connections, either of the specified database or of all databases, have cleared the `close_needed` state (see [the section called “SHOW SERVERS”](#)). This can be called after a `RECONNECT` or `RELOAD` to wait until the respective configuration change has been fully activated, for example in switchover scripts.

#### **Other Commands**

##### **SET *key* = *arg***

Changes a configuration setting (see also [the section called “SHOW CONFIG”](#)). For example:

```
SET log_connections = 1;
SET server_check_query = 'select 2';
```

(Note that this command is run on the pgbouncer admin console and sets pgbouncer settings. A `SET` command run on another database will be passed to the Postgres Pro backend like any other SQL command.)

#### **Signals**

`SIGHUP`

Reload config. Same as issuing the command `RELOAD` on the console.

SIGINT

Safe shutdown. Same as issuing `PAUSE` and `SHUTDOWN` on the console.

SIGTERM

Immediate shutdown. Same as issuing `SHUTDOWN` on the console.

SIGUSR1

Same as issuing `PAUSE` on the console.

SIGUSR2

Same as issuing `RESUME` on the console.

## Libevent Settings

From the libevent documentation:

It is possible to disable support for `epoll`, `kqueue`, `devpoll`, `poll`, or `select` by setting the environment variable `EVENT_NOEPOLL`, `EVENT_NOKQUEUE`, `EVENT_NODEVPOLL`, `EVENT_NOPOLL` or `EVENT_NOSELECT`, respectively.

By setting the environment variable `EVENT_SHOW_METHOD`, libevent displays the kernel notification method that it uses.

## pgbouncer.ini Configuration File

The configuration file is in the `.ini` format. Section names are between `[` and `]`. Lines starting with `;` or `#` are taken as comments and ignored. The characters `;` and `#` are not recognized as special when they appear later in the line.

### Generic Settings

`logfile`

Specifies the log file. For daemonization (`-d`), either this or `syslog` has to be set. The log file is kept open, so after rotation, `kill -HUP` or on console `RELOAD;` should be done. On Windows, the service must be stopped and started.

Note that setting `logfile` does not by itself turn off logging to `stderr`. Use the command-line option `-q` or `-d` for that.

Default: not set

`pidfile`

Specifies the PID file. Without `pidfile` set, daemonization (`-d`) is not allowed.

Default: not set

`listen_addr`

Specifies a list (comma-separated) of addresses where to listen for TCP connections. You may also use `*` meaning "listen on all addresses". When not set, only Unix socket connections are accepted.

Addresses can be specified numerically (IPv4/IPv6) or by name.

Default: not set

`listen_port`

Which port to listen on. Applies to both TCP and Unix sockets.

Default: 6432

#### `unix_socket_dir`

Specifies the location for Unix sockets. Applies to both the listening socket and server connections. If set to an empty string, Unix sockets are disabled. A value that starts with @ specifies that a Unix socket in the abstract namespace should be created (currently supported on Linux and Windows).

For online reboot (-R) to work, a Unix socket needs to be configured, and it needs to be in the file-system namespace.

Default: /tmp (empty on Windows)

#### `unix_socket_mode`

File system mode for Unix socket. Ignored for sockets in the abstract namespace. Not supported on Windows.

Default: 0777

#### `unix_socket_group`

Group name to use for Unix socket. Ignored for sockets in the abstract namespace. Not supported on Windows.

Default: not set

#### `user`

If set, specifies the Unix user to change to after startup. Works only if pgbouncer is started as root or if it's already running as the given user.

Not supported on Windows.

Default: not set

#### `pool_mode`

Specifies when a server connection can be reused by other clients.

##### `session`

Server is released back to pool after client disconnects. Default.

##### `transaction`

Server is released back to pool after transaction finishes.

##### `statement`

Server is released back to pool after query finishes. Transactions spanning multiple statements are disallowed in this mode.

#### `max_client_conn`

Maximum number of client connections allowed.

When this setting is increased, then the file descriptor limits in the operating system might also have to be increased. Note that the number of file descriptors potentially used is more than `max_client_conn`. If each user connects under its own username to the server, the theoretical maximum used is:

`max_client_conn + (max pool_size * total databases * total users)`

If a database user is specified in the connection string (all users connect under the same user name), the theoretical maximum is:

`max_client_conn + (max pool_size * total databases)`

The theoretical maximum should never be reached, unless somebody deliberately crafts a special load for it. Still, it means you should set the number of file descriptors to a safely high number.

Search for `ulimit` in your favorite shell man page. Note: `ulimit` does not apply in a Windows environment.

Default: 100

#### `default_pool_size`

How many server connections to allow per user/database pair. Can be overridden in the per-database configuration.

Default: 20

#### `min_pool_size`

Add more server connections to pool if below this number. Improves the behavior when the normal load suddenly comes back after a period of total inactivity. The value is effectively capped at the pool size.

Default: 0 (disabled)

#### `reserve_pool_size`

How many additional connections to allow to a pool (see `reserve_pool_timeout`). The 0 value disables this parameter.

Default: 0 (disabled)

#### `reserve_pool_timeout`

If a client has not been serviced in this time, `pgbouncer` enables use of additional connections from the reserve pool. The 0 value disables this parameter. [seconds]

Default: 5.0

#### `max_db_connections`

Do not allow more than this many server connections per database (regardless of user). This considers the `pgbouncer` database that the client has connected to, not the Postgres Pro database of the outgoing connection. This can also be set per database in the `[databases]` section.

Note that when you hit the limit, closing a client connection to one pool will not immediately allow a server connection to be established for another pool, because the server connection for the first pool is still open. Once the server connection closes (due to idle timeout), a new server connection will immediately be opened for the waiting pool.

Default: 0 (unlimited)

#### `max_user_connections`

Do not allow more than this many server connections per user (regardless of database). This considers the `pgbouncer` user that is associated with a pool, which is either the user specified for the server connection or in absence of that the user the client has connected as. This can also be set per user in the `[users]` section.

Note that when you hit the limit, closing a client connection to one pool will not immediately allow a server connection to be established for another pool, because the server connection for the first pool is still open. Once the server connection closes (due to idle timeout), a new server connection will immediately be opened for the waiting pool.

Default: 0 (unlimited)



### server\_round\_robin

By default, pgbouncer reuses server connections in LIFO (last-in, first-out) manner, so that few connections get the most load. This gives best performance if you have a single server serving a database. But if there is a round-robin system behind a database address (TCP, DNS, or host list), then it is better if pgbouncer also uses connections in that manner, thus achieving uniform load.

Default: 0

### track\_extra\_parameters

By default, pgbouncer tracks `client_encoding`, `datestyle`, `timezone`, `standard_conforming_strings` and `application_name` parameters per client. To allow other parameters to be tracked, they can be specified here, so that pgbouncer knows that they should be maintained in the client variable cache and restored in the server whenever the client becomes active.

If you need to specify multiple values, use a comma-separated list (e.g. `default_transaction_read-only, IntervalStyle`)

#### Note

Most parameters cannot be tracked this way. The only parameters that can be tracked are ones that Postgres Pro reports to the client. Postgres Pro has [an official list of parameters that it reports to the client](#). Postgres Pro extensions can change this list though, they can add parameters themselves that they also report, and they can start reporting already existing parameters that Postgres Pro does not report. Notably Citus 12.0+ causes PostgreSQL to also report `search_path`.

The `postgres` protocol allows specifying parameter settings, both directly as a parameter in the startup packet, or inside the [options startup packet](#). Parameters specified using both of these methods are supported by `track_extra_parameters`. However, it's not possible to include `options` itself in `track_extra_parameters`, only the parameters contained in `options`.

Default: `IntervalStyle`

### ignore\_startup\_parameters

By default, pgbouncer allows only parameters it can keep track of in startup packets: `client_encoding`, `datestyle`, `timezone` and `standard_conforming_strings`.

All other parameters will raise an error. To allow other parameters, they can be specified here, so that pgbouncer knows that they are handled by the admin and it can ignore them.

If you need to specify multiple values, use a comma-separated list (e.g. `options,extra_float_digits`).

The `postgres` protocol allows specifying parameter settings, both directly as a parameter in the startup packet, or inside the [options startup packet](#). Parameters specified using both of these methods are supported by `ignore_startup_parameters`. It's even possible to include `options` itself in `ignore_startup_parameters`, which results in any unknown parameters contained inside `options` to be ignored.

Default: empty

### peer\_id

The peer ID used to identify this pgbouncer process in a group of pgbouncer processes that are peered together. The `peer_id` value should be unique within a group of peered pgbouncer processes. When set to 0, pgbouncer peering is disabled. See also [\[peers\] Section](#) for more information. The maximum value that can be used for the `peer_id` is 16383.

Default: 0

`disable_pqexec`

Disable the Simple Query protocol (PQexec). Unlike the Extended Query protocol, Simple Query allows multiple queries in one packet, which allows some classes of SQL-injection attacks. Disabling it can improve security. Obviously, this means only clients that exclusively use the Extended Query protocol will stay working.

Default: 0

`application_name_add_host`

Add the client host address and port to the application name setting set on connection start. This helps in identifying the source of bad queries, etc. This logic applies only at the start of a connection. If `application_name` is later changed with `SET`, `pgbouncer` does not change it again.

Default: 0

`conffile`

Show location of current configuration file. Changing it will make `pgbouncer` use another configuration file for next `RELOAD` / `SIGHUP`.

Default: file from command line

`service_name`

Used on win32 service registration.

Default: `pgbouncer`

`job_name`

Alias for `service_name`.

`stats_period`

Sets how often the averages shown in various `SHOW` commands are updated and how often aggregated statistics are written to the log (but see `log_stats`). [seconds]

Default: 60

### Authentication Settings

`pgbouncer` handles its own client authentication and has its own database of users. These settings control this.

`auth_type`

How to authenticate users.

`cert`

The client must connect over TLS connection with a valid client certificate. The user name is then taken from the `CommonName` field from the certificate.

`md5`

Use MD5-based password check. This is the default authentication method. `auth_file` may contain both MD5-encrypted and plain-text passwords. If `md5` is configured and a user has a SCRAM secret, then SCRAM authentication is used automatically instead.

`scram-sha-256`

Use password check with SCRAM-SHA-256. `auth_file` has to contain SCRAM secrets or plain-text passwords. Note that SCRAM secrets can only be used for verifying the password of a client

but not for logging into a server. To be able to use SCRAM on server connections, use plain-text passwords.

#### plain

The clear-text password is sent over the wire. Deprecated.

#### trust

No authentication is done. The user name must still exist in `auth_file`.

#### any

Like the `trust` method, but the user name given is ignored. Requires that all databases are configured to log in as a specific user. Additionally, the console database allows any user to log in as admin.

#### hba

The actual authentication type is loaded from `auth_hba_file`. This allows different authentication methods for different access paths, for example: connections over Unix socket use `peer` authentication method, connections over TCP must use TLS.

#### pam

Pluggable Authentication Modules (PAM) method is used to authenticate users, `auth_file` is ignored. This method is not compatible with databases using the `auth_user` option. The service name reported to PAM is `pgbouncer`. `pam` is not supported in the HBA configuration file.

#### auth\_hba\_file

HBA configuration file to use when `auth_type` is `hba`.

Default: not set

#### auth\_file

The name of the file to load user names and passwords from. See [the section called “Authentication File Format”](#) for details.

Most authentication types (see `auth_type`) require that either `auth_file` or `auth_user` be set; otherwise there would be no users defined.

Default: not set

#### auth\_user

If `auth_user` is set, then any user not specified in `auth_file` will be queried through the `auth_query` query from `pg_shadow` in the database, using `auth_user`. The password of `auth_user` will be taken from `auth_file`. (If `auth_user` does not require a password, then it does not need to be defined in `auth_file`.)

Direct access to `pg_shadow` requires admin rights. It's preferable to use a non-superuser that calls a `SECURITY DEFINER` function instead.

Default: not set

#### auth\_query

Query to load user's password from database.

Direct access to `pg_shadow` requires admin rights. It's preferable to use a non-superuser that calls a `SECURITY DEFINER` function instead.

Note that the query is run inside the target database. So if a function is used, it needs to be installed into each database.

**Default:** `SELECT username, passwd FROM pg_shadow WHERE username=$1`

`auth_dbname`

Database name in the [\[databases\] section](#) to be used for authentication purposes. This option can be either global or overridden in the connection string if this parameter is specified.

### Log Settings

`syslog`

Toggles syslog on/off. On Windows, the event log is used instead.

**Default:** 0

`syslog_ident`

Under what name to send logs to syslog.

**Default:** `pgbouncer` (program name)

`syslog_facility`

Under what facility to send logs to syslog. Possibilities: `auth`, `authpriv`, `daemon`, `user`, `local0-7`.

**Default:** `daemon`

`log_connections`

Log successful logins.

**Default:** 1

`log_disconnections`

Log disconnections with reasons.

**Default:** 1

`log_pooler_errors`

Log error messages the pooler sends to clients.

**Default:** 1

`log_stats`

Write aggregated statistics into the log, every [stats\\_period](#). This can be disabled if external monitoring tools are used to grab the same data from `SHOW` commands.

**Default:** 1

`verbose`

Increase verbosity. Mirrors the `-v` switch on the command line. For example, using `-v -v` on the command line is the same as `verbose=2`.

**Default:** 0

### Console Access Control

`admin_users`

Comma-separated list of database users that are allowed to connect and run all commands on the console. Ignored when [auth\\_type](#) is `any`, in which case any user name is allowed in as admin.

**Default:** empty

#### stats\_users

Comma-separated list of database users that are allowed to connect and run read-only queries on the console. That means all `SHOW` commands except `SHOW FDS`.

Default: empty

### Connection Sanity Checks, Timeouts

#### server\_reset\_query

Query sent to server on connection release, before making it available to other clients. At that moment no transaction is in progress, so the value should not include `ABORT` or `ROLLBACK`.

The query is supposed to clean any changes made to the database session so that the next client gets the connection in a well-defined state. The default is `DISCARD ALL`, which cleans everything, but that leaves the next client no pre-cached state. It can be made lighter, e.g. `DEALLOCATE ALL` to just drop prepared statements, if the application does not break when some state is kept around.

When transaction pooling is used, the `server_reset_query` is not used, because in that mode, clients must not use any session-based features, since each transaction ends up in a different connection and thus gets a different session state.

Default: `DISCARD ALL`

#### server\_reset\_query\_always

Whether `server_reset_query` should be run in all pooling modes. When this setting is off (default), the `server_reset_query` will be run only in pools that are in sessions-pooling mode. Connections in transaction-pooling mode should not have any need for a reset query.

This setting is for working around broken setups that run applications that use session features over a transaction-pooled pgbouncer. It changes non-deterministic breakage to deterministic breakage: clients always lose their state after each transaction.

Default: 0

#### server\_check\_delay

How long to keep released connections available for immediate re-use, without running `server_check_query` on it. If 0 then the query is always run.

Default: 30.0

#### server\_check\_query

Simple do-nothing query to check if the server connection is alive.

If an empty string, then sanity checking is disabled.

Default: `select 1`

#### server\_fast\_close

Disconnect a server in session pooling mode immediately or after the end of the current transaction if it is in `close_needed` mode (set by `RECONNECT`, `RELOAD` that changes connection settings, or DNS change), rather than waiting for the session end. In statement or transaction pooling mode, this has no effect since that is the default behavior there.

If because of this setting a server connection is closed before the end of the client session, the client connection is also closed. This ensures that the client notices that the session has been interrupted.

This setting makes connection configuration changes take effect sooner if session pooling and long-running sessions are used. The downside is that client sessions are liable to be interrupted by a

configuration change, so client applications will need logic to reconnect and reestablish session state. But note that no transactions will be lost, because running transactions are not interrupted, only idle sessions.

Default: 0

#### `server_lifetime`

The pooler will close an unused (not currently linked to any client connection) server connection that has been connected longer than this. Setting it to 0 means the connection is to be used only once, then closed. [seconds]

Default: 3600.0

#### `server_idle_timeout`

If a server connection has been idle more than this many seconds it will be closed. If 0 then timeout is disabled. [seconds]

Default: 600.0

#### `server_connect_timeout`

If connection and login don't finish in this amount of time, the connection will be closed. [seconds]

Default: 15.0

#### `server_login_retry`

If login to the server failed, because of failure to connect or from authentication, the pooler waits this much before retrying to connect. During the waiting interval, new clients trying to connect to the failing server will get an error immediately without another connection attempt. [seconds]

The purpose of this behavior is that clients don't unnecessarily queue up waiting for a server connection to become available if the server is not working. However, it also means that if a server is momentarily failing, for example during a restart or if the configuration was erroneous, then it will take at least this long until the pooler will consider connecting to it again. Planned events such as restarts should normally be managed using the `PAUSE` command to avoid this.

Default: 15.0

#### `client_login_timeout`

If a client connects but does not manage to log in in this amount of time, it will be disconnected. Mainly needed to avoid dead connections stalling `SUSPEND` and thus online restart. [seconds]

Default: 60.0

#### `autodb_idle_timeout`

If the automatically created (via `"*`") database pools have been unused this many seconds, they are freed. The negative aspect of that is that their statistics are also forgotten. [seconds]

Default: 3600.0

#### `dns_max_ttl`

How long DNS lookups can be cached. The actual DNS TTL is ignored. [seconds]

Default: 15.0

#### `dns_nxdomain_ttl`

How long DNS errors and `NXDOMAIN` DNS lookups can be cached. [seconds]

Default: 15.0

`dns_zone_check_period`

Period to check if a zone serial has changed.

pgbouncer can collect DNS zones from host names (everything after first dot) and then periodically check if the zone serial changes. If it notices changes, all host names under that zone are looked up again. If any host IP changes, its connections are invalidated.

Default: 0.0 (disabled)

`resolv_conf`

The location of a custom `resolv.conf` file. This is to allow specifying custom DNS servers and perhaps other name resolution options, independent of the global operating system configuration.

The parsing of the file is done by the DNS backend library, not pgbouncer, so see the library's documentation for details on allowed syntax and directives.

Default: empty (use operating system defaults)

## TLS Settings

`client_tls_sslmode`

TLS mode to use for connections from clients. TLS connections are disabled by default. When enabled, `client_tls_key_file` and `client_tls_cert_file` must be also configured to set up the key and certificate pgbouncer uses to accept client connections.

`disable`

Plain TCP. If client requests TLS, it's ignored. Default.

`allow`

If client requests TLS, it is used. If not, plain TCP is used. If the client presents a client certificate, it is not validated.

`prefer`

Same as `allow`.

`require`

The client must use TLS. If not, the client connection is rejected. If the client presents a client certificate, it is not validated.

`verify-ca`

Client must use TLS with valid client certificate.

`verify-full`

Same as `verify-ca`.

`client_tls_key_file`

Private key for pgbouncer to accept client connections.

Default: not set

`client_tls_cert_file`

Certificate for private key. Clients can validate it.

Default: not set

`client_tls_ca_file`

Root certificate file to validate client certificates.

Default: not set

`client_tls_protocols`

Which TLS protocol versions are allowed. Allowed values: `tlsv1.0`, `tlsv1.1`, `tlsv1.2`, `tlsv1.3`. Shortcuts: `all` (`tlsv1.0,tlsv1.1,tlsv1.2,tlsv1.3`), `secure` (`tlsv1.2,tlsv1.3`), `legacy` (`all`).

Default: `secure`

`client_tls_ciphers`

Allowed TLS ciphers, in OpenSSL syntax. Shortcuts: `default/secure`, `compat/legacy`, `insecure/all`, `normal`, `fast`.

Only connections using TLS version 1.2 and lower are affected. There is currently no setting that controls the cipher choices used by TLS version 1.3 connections.

Default: `fast`

`client_tls_ecdhcurve`

Elliptic Curve name to use for ECDH key exchanges.

Allowed values: `none` (DH is disabled), `auto` (256-bit ECDH), curve name.

Default: `auto`

`client_tls_dheparams`

DHE key exchange type.

Allowed values: `none` (DH is disabled), `auto` (2048-bit DH), `legacy` (1024-bit DH).

Default: `auto`

`server_tls_sslmode`

TLS mode to use for connections to Postgres Pro servers. The default mode is `prefer`.

`disable`

Plain TCP. TLS is not even requested from the server.

`prefer`

TLS connection is always requested first from Postgres Pro. If refused, the connection will be established over plain TCP. Server certificate is not validated. Default.

`require`

Connection must go over TLS. If server rejects it, plain TCP is not attempted. Server certificate is not validated.

`verify-ca`

Connection must go over TLS and server certificate must be valid according to `server_tls_ca_file`. Server host name is not checked against certificate.

`verify-full`

Connection must go over TLS and server certificate must be valid according to `server_tls_ca_file`. Server host name must match certificate information.



`server_tls_ca_file`

Root certificate file to validate Postgres Pro server certificates.

Default: not set

`server_tls_key_file`

Private key for pgbouncer to authenticate against Postgres Pro server.

Default: not set

`server_tls_cert_file`

Certificate for private key. Postgres Pro server can validate it.

Default: not set

`server_tls_protocols`

Which TLS protocol versions are allowed. Allowed values: `tlsv1.0`, `tlsv1.1`, `tlsv1.2`, `tlsv1.3`. Shortcuts: `all` (`tlsv1.0,tlsv1.1,tlsv1.2,tlsv1.3`), `secure` (`tlsv1.2,tlsv1.3`), `legacy` (`all`).

Default: `secure`

`server_tls_ciphers`

Allowed TLS ciphers, in OpenSSL syntax. Shortcuts: `default/secure`, `compat/legacy`, `insecure/all`, `normal`, `fast`.

Only connections using TLS version 1.2 and lower are affected. There is currently no setting that controls the cipher choices used by TLS version 1.3 connections.

Default: `fast`

### **Dangerous Timeouts**

Setting the following timeouts can cause unexpected errors.

`query_timeout`

Queries running longer than that are canceled. This should be used only with a slightly smaller server-side `statement_timeout`, to apply only for network problems. [seconds]

Default: 0.0 (disabled)

`query_wait_timeout`

Maximum time queries are allowed to spend waiting for execution. If the query is not assigned to a server during that time, the client is disconnected. The 0 value disables this parameter. If this is disabled, clients will be queued indefinitely. [seconds]

This setting is used to prevent unresponsive servers from grabbing up connections. It also helps when the server is down or rejects connections for any reason.

Default: 120.0

`cancel_wait_timeout`

Maximum time cancellation requests are allowed to spend waiting for execution. If the cancel request is not assigned to a server during that time, the client is disconnected. 0 disables. If this is disabled, cancel requests will be queued indefinitely. [seconds]

This setting is used to prevent a client locking up when a cancel cannot be forwarded due to the server being down.

Default: 10.0

#### `client_idle_timeout`

Client connections idling longer than this many seconds are closed. This should be larger than the client-side connection lifetime settings, and only used for network problems. [seconds]

Default: 0.0 (disabled)

#### `idle_transaction_timeout`

If a client has been in the “idle in transaction” state longer, it will be disconnected. [seconds]

Default: 0.0 (disabled)

#### `suspend_timeout`

How long to wait for buffer flush during `SUSPEND` or reboot (`-R`). A connection is dropped if the flush does not succeed. [seconds]

Default: 10

### Low-Level Network Settings

#### `pkt_buf`

Internal buffer size for packets. Affects size of TCP packets sent and general memory usage. Actual libpq packets can be larger than this, so no need to set it large.

Default: 4096

#### `max_packet_size`

Maximum size for Postgres Pro packets that pgbouncer allows through. One packet is either one query or one result set row. The full result set can be larger.

Default: 2147483647

#### `listen_backlog`

The value of the `backlog` argument for `listen()`. Determines how many new unanswered connection attempts are kept in the queue. When the queue is full, further new connections are dropped.

Default: 128

#### `sbuf_loopcnt`

How many times to process data on one connection, before proceeding. Without this limit, one connection with a big result set can stall pgbouncer for a long time. One loop processes one `pkt_buf` amount of data. 0 means no limit.

Default: 5

#### `so_reuseport`

Specifies whether to set the socket option `SO_REUSEPORT` on TCP listening sockets. On some operating systems, this allows running multiple pgbouncer instances on the same host listening on the same port and having the kernel distribute the connections automatically. This option is a way to get pgbouncer to use more CPU cores. (pgbouncer is single-threaded and uses one CPU core per instance.)

This setting has the desired effect on Linux. On systems that don't support the socket option at all, turning this setting on will result in an error.

Each pgbouncer instance on the same host needs different settings for at least `unix_socket_dir` and `pidfile`, as well as `logfile` if that is used. Also note that if you make use of this option, you can no longer connect to a specific pgbouncer instance via TCP/IP, which might have implications for monitoring and metrics collection.

To make sure query cancellations keep working, you should set up pgbouncer peering between the different pgbouncer processes. For details see the [peer\\_id configuration option](#) and the [\[peers\] configuration section](#). There's also an example that uses peering and `so_reuseport` in the [Examples section](#).

Default: 0

`tcp_defer_accept`

Sets the `TCP_DEFER_ACCEPT` socket option; see `man 7 tcp` for details. (This is a Boolean option: 1 means enabled. The actual value set if enabled is currently hardcoded to 45 seconds.)

This is currently only supported on Linux.

Default: 1 on Linux, otherwise 0

`tcp_socket_buffer`

Default: not set

`tcp_keepalive`

Turns on basic keepalive with OS defaults.

On Linux, the system defaults are **`tcp_keepidle=7200`**, **`tcp_keepintvl=75`**, **`tcp_keepcnt=9`**. They are probably similar on other operating systems.

Default: 1

`tcp_keepcnt`

Default: not set

`tcp_keepidle`

Default: not set

`tcp_keepintvl`

Default: not set

`tcp_user_timeout`

Sets the `TCP_USER_TIMEOUT` socket option. This specifies the maximum amount of time in milliseconds that transmitted data may remain unacknowledged before the TCP connection is forcibly closed. If set to 0, then operating system's default is used.

This is currently only supported on Linux.

Default: 0

### Section [databases]

The section `[databases]` defines the names of the databases that clients of pgbouncer can connect to and specifies where those connections will be routed. The section contains `key=value` lines like

```
dbname = connection string
```

where the key will be taken as a database name and the value as a connection string, consisting of `key=value` pairs of connection parameters, described below (similar to `libpq`, but the actual `libpq` is not used and the set of available features is different).

Example:

```
foodb = host=host1.example.com port=5432
```

```
bardb = host=localhost dbname=bazdb
```

The database name can contain characters `_0-9A-Za-z` without quoting. Names that contain other chars need to be quoted with standard SQL ident quoting: double quotes where `"` is taken as single quote.

The database name `pgbouncer` is reserved for the admin console and cannot be used as a key here.

\* acts as fallback database: if the exact name does not exist, its value is taken as connection string for the requested database. For example, if there is the following entry (and no other overriding entries):

```
* = host=foo
```

In this case, a connection to `pgbouncer` specifying a database `bar` will effectively behave as if the following entry exists (taking advantage of the default for `dbname` being the client-side database name):

```
bar = host=foo dbname=bar
```

Such automatically created database entries are cleaned up if they stay idle longer than the time specified by the `autodb_idle_timeout` parameter.

`dbname`

Destination database name.

Default: same as client-side database name

`host`

Host name or IP address to connect to. Host names are resolved at connection time, the result is cached per `dns_max_ttl` parameter. When a host name's resolution changes, existing server connections are automatically closed when they are released (according to the pooling mode), and new server connections immediately use the new resolution. If DNS returns several results, they are used in a round-robin manner.

If the value begins with `/`, then a Unix socket in the file-system namespace is used. If the value begins with `@`, then a Unix socket in the abstract namespace is used.

A comma-separated list of host names or addresses can be specified. In that case, connections are made in a round-robin manner. (If a host list contains host names that in turn resolve via DNS to multiple addresses, the round-robin systems operate independently. This is an implementation dependency that is subject to change.) Note that in a list, all hosts must be available at all times: there are no mechanisms to skip unreachable hosts or to select only available hosts from a list or similar. (This is different from what a host list in `libpq` means.) Also note that this only affects how the destinations of new connections are chosen. See also the setting `server_round_robin` for how clients are assigned to already established server connections.

Examples:

```
host=localhost
host=127.0.0.1
host=2001:0db8:85a3:0000:0000:8a2e:0370:7334
host=/var/run/postgresql
host=192.168.0.1,192.168.0.2,192.168.0.3
```

Default: not set, meaning to use a Unix socket

`port`

Default: 5432

`user`

If `user=` is set, all connections to the destination database will be done with the specified user, meaning that there will be only one pool for this database.

Otherwise pgbouncer logs into the destination database with the client user name, meaning that there will be one pool per user.

password

If no password is specified here, the password from the `auth_file` or `auth_query` will be used.

auth\_user

Override of the global `auth_user` setting, if specified.

pool\_size

Set the maximum size of pools for this database. If not set, the `default_pool_size` is used.

min\_pool\_size

Set the minimum pool size for this database. If not set, the global `min_pool_size` is used.

reserve\_pool

Set additional connections for this database. If not set, `reserve_pool_size` is used.

connect\_query

Query to be executed after a connection is established, but before allowing the connection to be used by any clients. If the query raises errors, they are logged but ignored otherwise.

pool\_mode

Set the pool mode specific to this database. If not set, the default `pool_mode` is used.

max\_db\_connections

Configure a database-wide maximum (i.e. all pools within the database will not have more than this many server connections).

client\_encoding

Ask specific `client_encoding` from server.

datestyle

Ask specific `datestyle` from server.

timezone

Ask specific `timezone` from server.

### Section [users]

This section contains `key=value` lines like

```
user1 = settings
```

where the key will be taken as a user name and the value as a list of configuration settings specific for this user.

Example:

```
user1 = pool_mode=session
```

Only a few settings are available here.

pool\_mode

Set the pool mode to be used for all connections from this user. If not set, the database or default `pool_mode` is used.

`max_user_connections`

Configure a maximum for the user (i.e. all pools with the user will not have more than this many server connections).

### Section [peers]

This section defines the peers that pgbouncer can forward cancellation requests to and where those cancellation requests will be routed.

pgbouncer processes can be peered together in a group by defining a `peer_id` value and a [peers] section in the configs of all the pgbouncer processes. These pgbouncer processes can then forward cancellation requests to the process that it originated from. This is needed to make cancellations work when multiple pgbouncer processes (possibly on different servers) are behind the same TCP load balancer. Cancellation requests are sent over different TCP connections than the query they are cancelling, so a TCP load balancer might send the cancellation request connection to a different process than the one that it was meant for. By peering them these cancellation requests eventually end up at the right process.

The section contains key=value lines like

```
peer_id = connection string
```

where the key will be taken as a `peer_id` and the value as a connection string, consisting of key=value pairs of connection parameters, described below (similar to libpq, but the actual libpq is not used and the set of available features is different).

Example:

```
1 = host=host1.example.com
2 = host=/tmp/pgbouncer-2 port=5555
```

### Note

For peering to work, the `peer_id` of each pgbouncer process in the group must be unique within the peered group. And the [peers] section should contain entries for each of those peer IDs. An example can be found in the [Examples](#) section. It is allowed, but not necessary, for the [peers] section to contain the `peer_id` of the pgbouncer that the config is for. Such an entry will be ignored, but it is allowed to make config management easier. Because it allows using the exact same [peers] section for multiple configs.

`host`

Host name or IP address to connect to. Host names are resolved at connection time, the result is cached per `dns_max_ttl` parameter. If DNS returns several results, they are used in a round-robin manner. But in general it's not recommended to use a hostname that resolves to multiple IPs, because then the cancel request might still be forwarded to the wrong node and it would need to be forwarded again (which is only allowed up to three times).

If the value begins with `/`, then a Unix socket in the file-system namespace is used. If the value begins with `@`, then a Unix socket in the abstract namespace is used.

Examples:

```
host=localhost
host=127.0.0.1
host=2001:0db8:85a3:0000:0000:8a2e:0370:7334
host=/var/run/pgbouncer-1
```

`port`

Default: 6432

`pool_size`

If not set, the `default_pool_size` is used.

### Include Directive

The pgbouncer configuration file can contain include directives, which specify another configuration file to read and process. This allows splitting the configuration file into physically separate parts. The include directives look like this:

```
%include filename
```

If the filename is not an absolute path, it is taken as relative to the current working directory.

### Authentication File Format

This section describes the format of the file specified by the `auth_file` setting. It is a text file in the following format:

```
"username1" "password" ...
"username2" "md5abcdef012342345" ...
"username2" "SCRAM-SHA-256$iterations:salt$storedkey:serverkey"
```

There should be at least two fields, surrounded by double quotes. The first field is the user name and the second is either a plain-text, a MD5-hashed password, or a SCRAM secret. pgbouncer ignores the rest of the line. Double quotes in a field value can be escaped by writing two double quotes.

Postgres Pro MD5-hashed password format:

```
"md5" + md5(password + username)
```

```
So user admin with password 1234 will have MD5-hashed password
md545f2603610af569b6155c45067268c6b.
```

Postgres Pro SCRAM secret format:

```
SCRAM-SHA-256$iterations:salt$storedkey:serverkey
```

The passwords or secrets stored in the authentication file serve two purposes. First, they are used to verify the passwords of incoming client connections, if a password-based authentication method is configured. Second, they are used as the passwords for outgoing connections to the backend server, if the backend server requires password-based authentication (unless the password is specified directly in the database's connection string). The latter works if the password is stored in plain text or MD5-hashed. SCRAM secrets can only be used for logging into a server if the client authentication also uses SCRAM, the pgbouncer database definition does not specify a user name, and the SCRAM secrets are identical in pgbouncer and the Postgres Pro server (same salt and iterations, not merely the same password). This is due to an inherent security property of SCRAM: the stored SCRAM secret cannot by itself be used for deriving login credentials.

The authentication file can be written by hand, but it's also useful to generate it from some other list of users and passwords. See `./etc/mkauth.py` for a sample script to generate the authentication file from the `pg_shadow` system table.

Alternatively, use `auth_query` instead of `auth_file` to avoid having to maintain a separate authentication file.

### HBA File Format

The location of the HBA file is specified by the setting `auth_hba_file`. It is only used if `auth_type` is set to `hba`.

The file follows the format of the Postgres Pro `pg_hba.conf` file described in [Section 20.1](#).

- Supported record types: `local`, `host`, `hostssl`, `hostnossl`.

- **Database field:** Supports `all`, `sameuser`, `@file`, multiple names. Not supported: `replication`, `samerole`, `samegroup`.
- **User name field:** Supports `all`, `@file`, multiple names. Not supported: `+groupname`.
- **Address field:** Supports `IPv4`, `IPv6`. Not supported: DNS names, domain prefixes.
- **Auth-method field:** Only methods supported by `pgbouncer`'s `auth_type` are supported, plus `peer` and `reject`, but except `any` and `pam`, which only work globally. User name map (`map=`) parameter is not supported.

### Examples

Small example configuration:

```
[databases]
template1 = host=localhost dbname=template1 auth_user=someuser

[pgbouncer]
pool_mode = session
listen_port = 6432
listen_addr = localhost
auth_type = md5
auth_file = users.txt
logfile = pgbouncer.log
pidfile = pgbouncer.pid
admin_users = someuser
stats_users = stat_collector
```

Database examples:

```
[databases]

; foodb over Unix socket
foodb =

; redirect bardb to bazdb on localhost
bardb = host=localhost dbname=bazdb

; access to destination database will go with single user
forcedb = host=localhost port=300 user=baz password=foo client_encoding=UNICODE
datestyle=ISO
```

Example of a secure function for `auth_query`:

```
CREATE OR REPLACE FUNCTION pgbouncer.user_lookup(in i_username text, out uname text,
out phash text)
RETURNS record AS $$
BEGIN
 SELECT username, passwd FROM pg_catalog.pg_shadow
 WHERE username = i_username INTO uname, phash;
 RETURN;
END; $$ LANGUAGE plpgsql SECURITY DEFINER;
REVOKE ALL ON FUNCTION pgbouncer.user_lookup(text) FROM public, pgbouncer;
GRANT EXECUTE ON FUNCTION pgbouncer.user_lookup(text) TO pgbouncer;
```

Example configs for 2 peered `pgbouncer` processes to create a multi-core `pgbouncer` setup using [so\\_reuseport](#).

The config for the first process:

```
[databases]
postgres = host=localhost dbname=postgres
```



```
[peers]
1 = host=/tmp/pgbouncer1
2 = host=/tmp/pgbouncer2

[pgbouncer]
listen_addr=127.0.0.1
auth_file=auth_file.conf
so_reuseport=1
unix_socket_dir=/tmp/pgbouncer1
peer_id=1
```

**The config for the second process:**

```
[databases]
postgres = host=localhost dbname=postgres

[peers]
1 = host=/tmp/pgbouncer1
2 = host=/tmp/pgbouncer2

[pgbouncer]
listen_addr=127.0.0.1
auth_file=auth_file.conf
so_reuseport=1
; only unix_socket_dir and peer_id are different
unix_socket_dir=/tmp/pgbouncer2
peer_id=2
```

---

# Appendix I. External Projects

Postgres Pro is a complex software project, and managing the project is difficult. We have found that many enhancements to Postgres Pro can be more efficiently developed separately from the core project.

## I.1. Client Interfaces

There are only two client interfaces included in the base Postgres Pro Shardman distribution:

- [libpq](#) is included because it is the primary C language interface, and because many other client interfaces are built on top of it.
- [ECPG](#) is included because it depends on the server-side SQL grammar, and is therefore sensitive to changes in Postgres Pro itself.

All other language interfaces are external projects and are distributed separately. A [list of language interfaces](#) is maintained on the PostgreSQL wiki. Note that some of these packages are not released under the same license as Postgres Pro. For more information on each language interface, including licensing terms, refer to its website and documentation.

[https://wiki.postgresql.org/wiki/List\\_of\\_drivers](https://wiki.postgresql.org/wiki/List_of_drivers)

## I.2. Administration Tools

There are several administration tools available for Postgres Pro. One of them is PPEM, an integrated administration panel for Postgres Pro Shardman management. It provides centralized privilege-based access to a unified management and monitoring interface with embedded command line. There are other solutions, e.g. an open-source platform [pgAdmin](#).

## I.3. Procedural Languages

Postgres Pro includes several procedural languages with the base distribution: [PL/pgSQL](#), [PL/Tcl](#), [PL/Perl](#), and [PL/Python](#).

In addition, there are a number of procedural languages that are developed and maintained outside the core Postgres Pro distribution. A list of [procedural languages](#) is maintained on the PostgreSQL wiki. Note that some of these projects are not released under the same license as Postgres Pro. For more information on each procedural language, including licensing information, refer to its website and documentation.

[https://wiki.postgresql.org/wiki/PL\\_Matrix](https://wiki.postgresql.org/wiki/PL_Matrix)

## I.4. Extensions

Postgres Pro is designed to be easily extensible. For this reason, extensions loaded into the database can function just like features that are built in. The extensions shipped with Postgres Pro are described in [Appendix F](#). Other extensions are developed independently, like [PostGIS](#). Even Postgres Pro replication solutions can be developed externally. For example, [Slony-I](#) is a popular primary/standby replication solution that is developed independently from the core project. Some extensions are made compatible with Postgres Pro Shardman. Consult the Postgres Pro Shardman support team if you want to gain access to such extensions.

---

# Appendix J. Postgres Pro Limits

[Table J.1](#) describes various hard limits of Postgres Pro. However, practical limits, such as performance limitations or available disk space may apply before absolute hard limits are reached.

**Table J.1. Postgres Pro Limitations**

| Item                    | Upper Limit                                                           | Comment                                                                |
|-------------------------|-----------------------------------------------------------------------|------------------------------------------------------------------------|
| database size           | unlimited                                                             |                                                                        |
| number of databases     | 4,294,950,911                                                         |                                                                        |
| relations per database  | 1,431,650,303                                                         |                                                                        |
| relation size           | 32 TB                                                                 | with the default <code>BLCKSZ</code> of 8192 bytes                     |
| rows per table          | limited by the number of tuples that can fit onto 4,294,967,295 pages |                                                                        |
| columns per table       | 1,600                                                                 | further limited by tuple size fitting on a single page; see note below |
| columns in a result set | 1,664                                                                 |                                                                        |
| field size              | 1 GB                                                                  |                                                                        |
| indexes per table       | unlimited                                                             | constrained by maximum relations per database                          |
| columns per index       | 32                                                                    | can be increased by recompiling Postgres Pro                           |
| partition keys          | 32                                                                    | can be increased by recompiling Postgres Pro                           |
| identifier length       | 63 bytes                                                              | can be increased by recompiling Postgres Pro                           |
| function arguments      | 100                                                                   | can be increased by recompiling Postgres Pro                           |
| query parameters        | 65,535                                                                |                                                                        |

The maximum number of columns for a table is further reduced as the tuple being stored must fit in a single 8192-byte heap page. For example, excluding the tuple header, a tuple made up of 1,600 `int` columns would consume 6400 bytes and could be stored in a heap page, but a tuple of 1,600 `bigint` columns would consume 12800 bytes and would therefore not fit inside a heap page. Variable-length fields of types such as `text`, `varchar`, and `char` can have their values stored out of line in the table's TOAST table when the values are large enough to require it. Only an 18-byte pointer must remain inside the tuple in the table's heap. For shorter length variable-length fields, either a 4-byte or 1-byte field header is used and the value is stored inside the heap tuple.

Columns that have been dropped from the table also contribute to the maximum column limit. Moreover, although the dropped column values for newly created tuples are internally marked as null in the tuple's null bitmap, the null bitmap also occupies space.

Each table can store a theoretical maximum of  $2^{32}$  out-of-line values; see [Section 66.2](#) for a detailed discussion of out-of-line storage. This limit arises from the use of a 32-bit OID to identify each such value. The practical limit is significantly less than the theoretical limit, because as the OID space fills up, finding an OID that is still free can become expensive, in turn slowing down INSERT/UPDATE statements. Typically, this is only an issue for tables containing many terabytes of data; partitioning is a possible workaround.

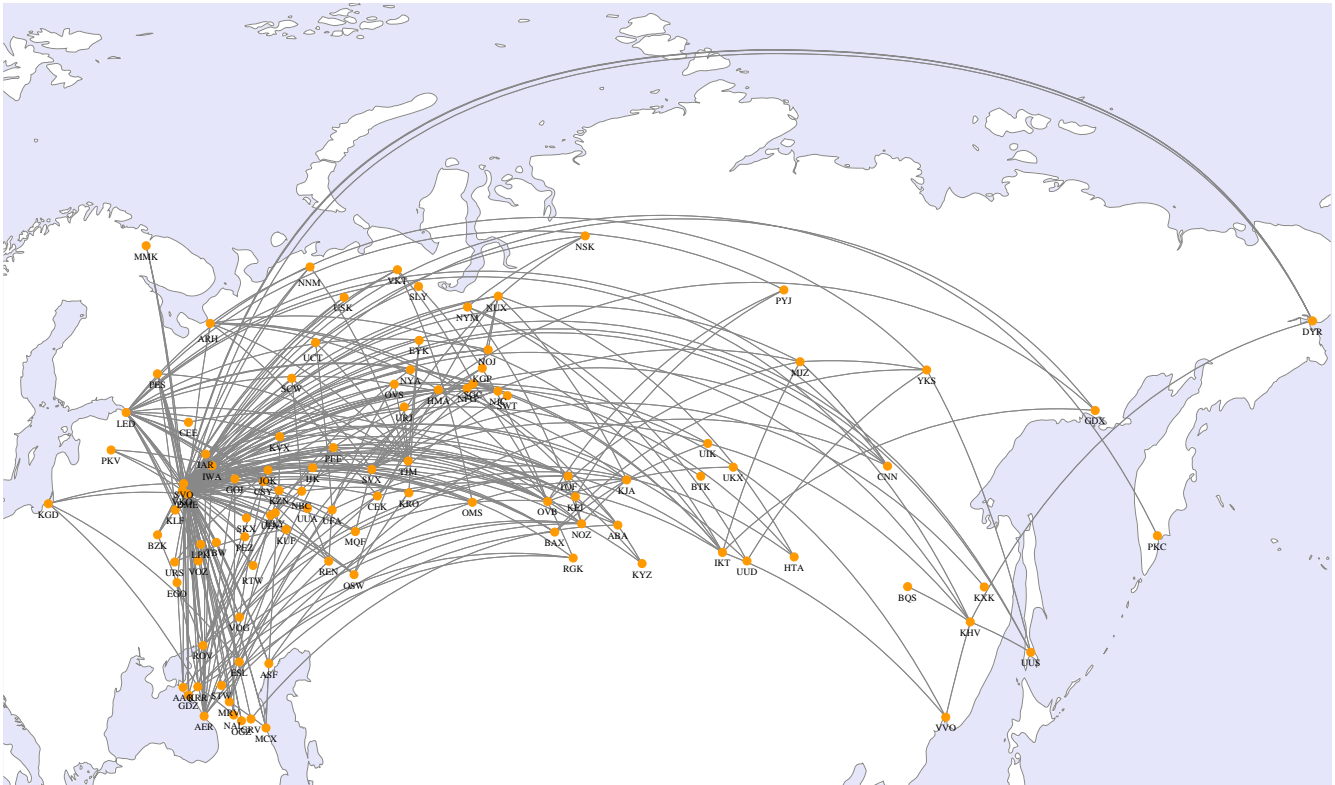
---

# Appendix K. Demo Database

## “Airlines”

This is an overview of a demo database for Postgres Pro. This appendix describes the database schema, which consists of eight tables and several views. The subject field of this database is airline flights in Russia. You can download the database from [our website](#). See [Section K.1](#) for details.

**Figure K.1. Airlines in Russia**



You can use this database for various purposes, such as:

- learning SQL language on your own
- preparing books, manuals, and courses on SQL
- showing Postgres Pro features in stories and articles

When developing this demo database, we pursued several goals:

- Database schema must be simple enough to be understood without extra explanations.
- At the same time, database schema must be complex enough to allow writing meaningful queries.
- The database must contain true-to-life data that will be interesting to work with.

This demo database is distributed under the [PostgreSQL license](#).

You can send us your feedback to [edu@postgrespro.ru](mailto:edu@postgrespro.ru).

## K.1. Installation

The demo database is available at [edu.postgrespro.com](http://edu.postgrespro.com) in three flavors, which differ only in the data size:

- [demo-small-en.zip](#) (21 MB) — flight data for one month (DB size is about 300 MB)

- [demo-medium-en.zip](#) (62 MB) — flight data for three months (DB size is about 700 MB)
- [demo-big-en.zip](#) (232 MB) — flight data for one year (DB size is about 2.5 GB)

The small database is good for writing queries, and it will not take up much disk space. The large database can help you understand the query behavior on large data volumes and consider query optimization.

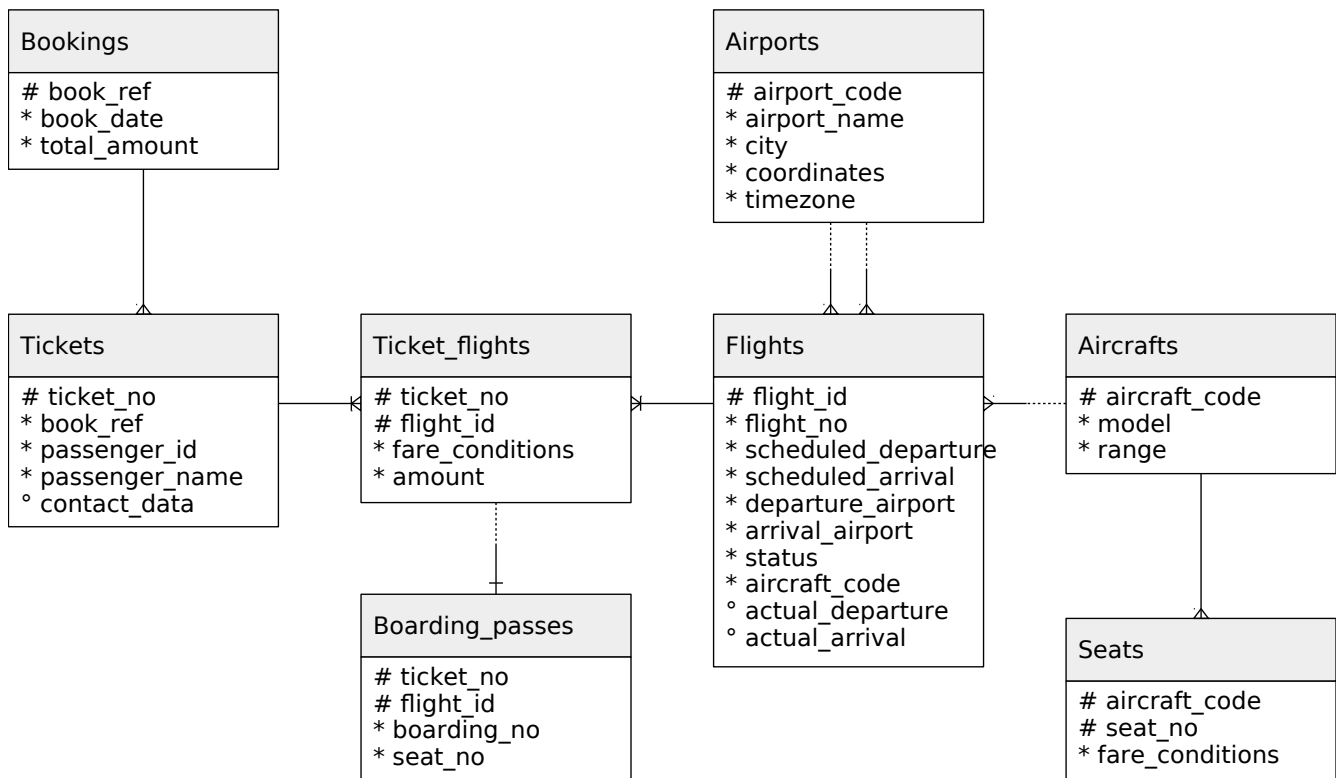
The files include an SQL script that creates the `demo` database and fills it with data (virtually, it is a backup copy created with the `pg_dump` utility). The owner of the `demo` database will be the DBMS user who runs the script. For example, to create the small database, run the script as the user `postgres` by means of `psql`:

```
psql -f demo_small_YYYYMMDD.sql -U postgres
```

Note that if the `demo` database already exists, it will be deleted and recreated!

## K.2. Schema Diagram

Figure K.2. Bookings Schema Diagram



## K.3. Schema Description

The main entity is a booking (`bookings`).

One booking can include several passengers, with a separate ticket (`tickets`) issued to each passenger. A ticket has a unique number and includes information about the passenger. As such, the passenger is not a separate entity. Both the passenger's name and identity document number can change over time, so it is impossible to uniquely identify all the tickets of a particular person; for simplicity, we can assume that all passengers are unique.

The ticket includes one or more flight segments (`ticket_flights`). Several flight segments can be included into a single ticket if there are no non-stop flights between the points of departure and destination (connecting flights), or if it is a round-trip ticket. Although there is no constraint in the schema, it is assumed that all tickets in the booking have the same flight segments.

Each flight (`flights`) goes from one airport (`airports`) to another. Flights with the same flight number have the same points of departure and destination, but differ in departure date.

At flight check-in, the passenger is issued a boarding pass (`boarding_passes`), where the seat number is specified. The passenger can check in for the flight only if this flight is included into the ticket. The flight-seat combination must be unique to avoid issuing two boarding passes for the same seat.

The number of seats (`seats`) in the aircraft and their distribution between different travel classes depends on the model of the aircraft (`aircrafts`) performing the flight. It is assumed that every aircraft model has only one cabin configuration. Database schema does not check that seat numbers in boarding passes have the corresponding seats in the aircraft (such verification can be done using table triggers, or at the application level).

## K.4. Schema Objects

### K.4.1. List of Relations

| Name                         | Type  | Small | Medium | Big    | Description             |
|------------------------------|-------|-------|--------|--------|-------------------------|
| <code>aircrafts</code>       | view  |       |        |        | Aircraft                |
| <code>aircrafts_data</code>  | table | 16 kB | 16 kB  | 16 kB  | Aircraft (translations) |
| <code>airports</code>        | view  |       |        |        | Airports                |
| <code>airports_data</code>   | table | 56 kB | 56 kB  | 56 kB  | Airports (translations) |
| <code>boarding_passes</code> | table | 31 MB | 102 MB | 427 MB | Boarding passes         |
| <code>bookings</code>        | table | 13 MB | 30 MB  | 105 MB | Bookings                |
| <code>flights</code>         | table | 3 MB  | 6 MB   | 19 MB  | Flights                 |
| <code>flights_v</code>       | view  |       |        |        | Flights                 |
| <code>routes</code>          | view  |       |        |        | Routes                  |
| <code>seats</code>           | table | 88 kB | 88 kB  | 88 kB  | Seats                   |
| <code>ticket_flights</code>  | table | 64 MB | 145 MB | 516 MB | Flight segments         |
| <code>tickets</code>         | table | 47 MB | 107 MB | 381 MB | Tickets                 |

### K.4.2. View `bookings.aircrafts`

Each aircraft model is identified by its three-digit code (`aircraft_code`). The view also includes the name of the aircraft model (`model`) and the maximal flying distance, in kilometers (`range`).

The value of the `model` field is selected according to the chosen language. See [Section K.4.15](#) for details.

| Column                     | Type                 | Modifiers             | Description                 |
|----------------------------|----------------------|-----------------------|-----------------------------|
| <code>aircraft_code</code> | <code>char(3)</code> | <code>not null</code> | Aircraft code, IATA         |
| <code>model</code>         | <code>text</code>    | <code>not null</code> | Aircraft model              |
| <code>range</code>         | <code>integer</code> | <code>not null</code> | Maximal flying distance, km |

View definition:

```
SELECT ml.aircraft_code,
 ml.model ->> lang() AS model,
 ml.range
FROM aircrafts_data ml;
```

### K.4.3. Table `bookings.aircrafts_data`

This is the base table for the `aircrafts` view. The `model` field of this table contains translations of aircraft models to different languages, in the JSONB format. In most cases, this table is not supposed to be used directly.

| Column | Type | Modifiers | Description |
|--------|------|-----------|-------------|
|--------|------|-----------|-------------|

```

-----+-----+-----+-----
aircraft_code | char(3) | not null | Aircraft code, IATA
model | jsonb | not null | Aircraft model
range | integer | not null | Maximal flying distance, km
Indexes:
 PRIMARY KEY, btree (aircraft_code)
Check constraints:
 CHECK (range > 0)
Referenced by:
 TABLE "flights" FOREIGN KEY (aircraft_code)
 REFERENCES aircrafts_data(aircraft_code)
 TABLE "seats" FOREIGN KEY (aircraft_code)
 REFERENCES aircrafts_data(aircraft_code) ON DELETE CASCADE

```

#### K.4.4. View `bookings.airports`

An airport is identified by a three-letter code (`airport_code`) and has a name (`airport_name`).

There is no separate entity for the city, but there is a city name (`city`) to identify the airports of the same city. The view also includes coordinates (`coordinates`) and the time zone (`timezone`).

The values of the `airport_name` and `city` fields are selected according to the chosen language. See [Section K.4.15](#) for details.

| Column                    | Type                 | Modifiers             | Description                                  |
|---------------------------|----------------------|-----------------------|----------------------------------------------|
| <code>airport_code</code> | <code>char(3)</code> | <code>not null</code> | Airport code                                 |
| <code>airport_name</code> | <code>text</code>    | <code>not null</code> | Airport name                                 |
| <code>city</code>         | <code>text</code>    | <code>not null</code> | City                                         |
| <code>coordinates</code>  | <code>point</code>   | <code>not null</code> | Airport coordinates (longitude and latitude) |
| <code>timezone</code>     | <code>text</code>    | <code>not null</code> | Airport time zone                            |

```

View definition:
SELECT ml.airport_code,
 ml.airport_name ->> lang() AS airport_name,
 ml.city ->> lang() AS city,
 ml.coordinates,
 ml.timezone
FROM airports_data ml;

```

#### K.4.5. Table `bookings.airports_data`

This is the base table for the `airports` view. This table contains translations of `airport_name` and `city` values to different languages, in the JSONB format. In most cases, this table is not supposed to be used directly.

| Column                    | Type                 | Modifiers             | Description                                  |
|---------------------------|----------------------|-----------------------|----------------------------------------------|
| <code>airport_code</code> | <code>char(3)</code> | <code>not null</code> | Airport code                                 |
| <code>airport_name</code> | <code>jsonb</code>   | <code>not null</code> | Airport name                                 |
| <code>city</code>         | <code>jsonb</code>   | <code>not null</code> | City                                         |
| <code>coordinates</code>  | <code>point</code>   | <code>not null</code> | Airport coordinates (longitude and latitude) |
| <code>timezone</code>     | <code>text</code>    | <code>not null</code> | Airport time zone                            |

```

Indexes:
 PRIMARY KEY, btree (airport_code)
Referenced by:
 TABLE "flights" FOREIGN KEY (arrival_airport)
 REFERENCES airports_data(airport_code)

```

```
TABLE "flights" FOREIGN KEY (departure_airport)
REFERENCES airports_data(airport_code)
```

#### K.4.6. Table `bookings.boarding_passes`

At the time of check-in, which opens twenty-four hours before the scheduled departure, the passenger is issued a boarding pass. Like the flight segment, the boarding pass is identified by the ticket number and the flight number.

Boarding passes are assigned sequential numbers (`boarding_no`), in the order of check-ins for the flight (this number is unique only within the context of a particular flight). The boarding pass specifies the seat number (`seat_no`).

| Column                   | Type                    | Modifiers             | Description          |
|--------------------------|-------------------------|-----------------------|----------------------|
| <code>ticket_no</code>   | <code>char(13)</code>   | <code>not null</code> | Ticket number        |
| <code>flight_id</code>   | <code>integer</code>    | <code>not null</code> | Flight ID            |
| <code>boarding_no</code> | <code>integer</code>    | <code>not null</code> | Boarding pass number |
| <code>seat_no</code>     | <code>varchar(4)</code> | <code>not null</code> | Seat number          |

Indexes:

```
PRIMARY KEY, btree (ticket_no, flight_id)
UNIQUE CONSTRAINT, btree (flight_id, boarding_no)
UNIQUE CONSTRAINT, btree (flight_id, seat_no)
```

Foreign-key constraints:

```
FOREIGN KEY (ticket_no, flight_id)
REFERENCES ticket_flights(ticket_no, flight_id)
```

#### K.4.7. Table `bookings.bookings`

Passengers book tickets for themselves, and, possibly, for several other passengers, in advance (`book_date`, not earlier than one month before the flight). The booking is identified by its number (`book_ref`, a six-position combination of letters and digits).

The `total_amount` field stores the total cost of all tickets included into the booking, for all passengers.

| Column                    | Type                       | Modifiers             | Description        |
|---------------------------|----------------------------|-----------------------|--------------------|
| <code>book_ref</code>     | <code>char(6)</code>       | <code>not null</code> | Booking number     |
| <code>book_date</code>    | <code>timestampz</code>    | <code>not null</code> | Booking date       |
| <code>total_amount</code> | <code>numeric(10,2)</code> | <code>not null</code> | Total booking cost |

Indexes:

```
PRIMARY KEY, btree (book_ref)
```

Referenced by:

```
TABLE "tickets" FOREIGN KEY (book_ref) REFERENCES bookings(book_ref)
```

#### K.4.8. Table `bookings.flights`

The natural key of the `bookings.flights` table consists of two fields — `flight_no` and `scheduled_departure`. To make foreign keys for this table more compact, a surrogate key is used as the primary key (`flight_id`).

A flight always connects two points — the airport of departure (`departure_airport`) and arrival (`arrival_airport`). There is no such entity as a “connecting flight”: if there are no non-stop flights from one airport to another, the ticket simply includes several required flight segments.

Each flight has a scheduled date and time of departure (`scheduled_departure`) and arrival (`scheduled_arrival`). The actual departure time (`actual_departure`) and arrival time (`actual_arrival`) can



differ: the difference is usually not very big, but sometimes can be up to several hours if the flight is delayed.

**Flight status** (*status*) can take one of the following values:

Scheduled

The flight is available for booking. It happens one month before the planned departure date; before that time, there is no entry for this flight in the database.

On Time

The flight is open for check-in (in twenty-four hours before the scheduled departure) and is not delayed.

Delayed

The flight is open for check-in (in twenty-four hours before the scheduled departure) but is delayed.

Departed

The aircraft has already departed and is airborne.

Arrived

The aircraft has reached the point of destination.

Cancelled

The flight is canceled.

| Column              | Type        | Modifiers | Description              |
|---------------------|-------------|-----------|--------------------------|
| flight_id           | serial      | not null  | Flight ID                |
| flight_no           | char(6)     | not null  | Flight number            |
| scheduled_departure | timestampz  | not null  | Scheduled departure time |
| scheduled_arrival   | timestampz  | not null  | Scheduled arrival time   |
| departure_airport   | char(3)     | not null  | Airport of departure     |
| arrival_airport     | char(3)     | not null  | Airport of arrival       |
| status              | varchar(20) | not null  | Flight status            |
| aircraft_code       | char(3)     | not null  | Aircraft code, IATA      |
| actual_departure    | timestampz  |           | Actual departure time    |
| actual_arrival      | timestampz  |           | Actual arrival time      |

Indexes:

PRIMARY KEY, btree (flight\_id)

UNIQUE CONSTRAINT, btree (flight\_no, scheduled\_departure)

Check constraints:

CHECK (scheduled\_arrival > scheduled\_departure)

CHECK ((actual\_arrival IS NULL)

OR ((actual\_departure IS NOT NULL AND actual\_arrival IS NOT NULL)  
AND (actual\_arrival > actual\_departure)))

CHECK (status IN ('On Time', 'Delayed', 'Departed',  
'Arrived', 'Scheduled', 'Cancelled'))

Foreign-key constraints:

FOREIGN KEY (aircraft\_code)

REFERENCES aircrafts(aircraft\_code)

FOREIGN KEY (arrival\_airport)

REFERENCES airports(airport\_code)

FOREIGN KEY (departure\_airport)

REFERENCES airports(airport\_code)

Referenced by:

TABLE "ticket\_flights" FOREIGN KEY (flight\_id)

```
REFERENCES flights(flight_id)
```

### K.4.9. Table bookings.seats

Seats define the cabin configuration of each aircraft model. Each seat is defined by its number (`seat_no`) and has an assigned travel class (`fare_conditions`): Economy, Comfort or Business.

| Column          | Type        | Modifiers | Description         |
|-----------------|-------------|-----------|---------------------|
| aircraft_code   | char(3)     | not null  | Aircraft code, IATA |
| seat_no         | varchar(4)  | not null  | Seat number         |
| fare_conditions | varchar(10) | not null  | Travel class        |

Indexes:

```
PRIMARY KEY, btree (aircraft_code, seat_no)
```

Check constraints:

```
CHECK (fare_conditions IN ('Economy', 'Comfort', 'Business'))
```

Foreign-key constraints:

```
FOREIGN KEY (aircraft_code)
```

```
REFERENCES aircrafts(aircraft_code) ON DELETE CASCADE
```

### K.4.10. Table bookings.ticket\_flights

A flight segment connects a ticket with a flight and is identified by their numbers.

Each flight has its cost (`amount`) and travel class (`fare_conditions`).

| Column          | Type          | Modifiers | Description   |
|-----------------|---------------|-----------|---------------|
| ticket_no       | char(13)      | not null  | Ticket number |
| flight_id       | integer       | not null  | Flight ID     |
| fare_conditions | varchar(10)   | not null  | Travel class  |
| amount          | numeric(10,2) | not null  | Travel cost   |

Indexes:

```
PRIMARY KEY, btree (ticket_no, flight_id)
```

Check constraints:

```
CHECK (amount >= 0)
```

```
CHECK (fare_conditions IN ('Economy', 'Comfort', 'Business'))
```

Foreign-key constraints:

```
FOREIGN KEY (flight_id) REFERENCES flights(flight_id)
```

```
FOREIGN KEY (ticket_no) REFERENCES tickets(ticket_no)
```

Referenced by:

```
TABLE "boarding_passes" FOREIGN KEY (ticket_no, flight_id)
```

```
REFERENCES ticket_flights(ticket_no, flight_id)
```

### K.4.11. Table bookings.tickets

A ticket has a unique number (`ticket_no`) that consists of 13 digits.

The ticket includes a passenger ID (`passenger_id`) — the identity document number, — their first and last names (`passenger_name`), and contact information (`contact_data`).

Neither the passenger ID, nor the name is permanent (for example, one can change the last name or passport), so it is impossible to uniquely identify all tickets of a particular passenger.

| Column    | Type     | Modifiers | Description   |
|-----------|----------|-----------|---------------|
| ticket_no | char(13) | not null  | Ticket number |

|                |             |          |                               |
|----------------|-------------|----------|-------------------------------|
| book_ref       | char(6)     | not null | Booking number                |
| passenger_id   | varchar(20) | not null | Passenger ID                  |
| passenger_name | text        | not null | Passenger name                |
| contact_data   | jsonb       |          | Passenger contact information |

Indexes:

PRIMARY KEY, btree (ticket\_no)

Foreign-key constraints:

FOREIGN KEY (book\_ref) REFERENCES bookings(book\_ref)

Referenced by:

TABLE "ticket\_flights" FOREIGN KEY (ticket\_no) REFERENCES tickets(ticket\_no)

## K.4.12. View bookings.flights\_v

There is a `flights_v` view over the `flights` table that provides additional information:

- Details about the airport of departure — `departure_airport`, `departure_airport_name`, `departure_city`
- Details about the airport of arrival — `arrival_airport`, `arrival_airport_name`, `arrival_city`
- Local departure time — `scheduled_departure_local`, `actual_departure_local`
- Local arrival time — `scheduled_arrival_local`, `actual_arrival_local`
- Flight duration — `scheduled_duration`, `actual_duration`.

| Column                    | Type        | Description                                                       |
|---------------------------|-------------|-------------------------------------------------------------------|
| flight_id                 | integer     | Flight ID                                                         |
| flight_no                 | char(6)     | Flight number                                                     |
| scheduled_departure       | timestamp   | Scheduled departure time                                          |
| scheduled_departure_local | timestamp   | Scheduled departure time,<br>local time at the point of departure |
| scheduled_arrival         | timestamp   | Scheduled arrival time                                            |
| scheduled_arrival_local   | timestamp   | Scheduled arrival time,<br>local time at the point of destination |
| scheduled_duration        | interval    | Scheduled flight duration                                         |
| departure_airport         | char(3)     | Departure airport code                                            |
| departure_airport_name    | text        | Departure airport name                                            |
| departure_city            | text        | City of departure                                                 |
| arrival_airport           | char(3)     | Arrival airport code                                              |
| arrival_airport_name      | text        | Arrival airport name                                              |
| arrival_city              | text        | City of arrival                                                   |
| status                    | varchar(20) | Flight status                                                     |
| aircraft_code             | char(3)     | Aircraft code, IATA                                               |
| actual_departure          | timestamp   | Actual departure time                                             |
| actual_departure_local    | timestamp   | Actual departure time,<br>local time at the point of departure    |
| actual_arrival            | timestamp   | Actual arrival time                                               |
| actual_arrival_local      | timestamp   | Actual arrival time,<br>local time at the point of destination    |
| actual_duration           | interval    | Actual flight duration                                            |

## K.4.13. View bookings.routes

The `bookings.flights` table contains some redundancies, which you can use to single out route information (flight number, airports of departure and destination) that does not depend on the exact flight dates.

Such information is shown in the `routes` view.

| Column                              | Type                   | Description                                     |
|-------------------------------------|------------------------|-------------------------------------------------|
| <code>flight_no</code>              | <code>char(6)</code>   | Flight number                                   |
| <code>departure_airport</code>      | <code>char(3)</code>   | Departure airport code                          |
| <code>departure_airport_name</code> | <code>text</code>      | Departure airport name                          |
| <code>departure_city</code>         | <code>text</code>      | City of departure                               |
| <code>arrival_airport</code>        | <code>char(3)</code>   | Arrival airport code                            |
| <code>arrival_airport_name</code>   | <code>text</code>      | Arrival airport name                            |
| <code>arrival_city</code>           | <code>text</code>      | City of arrival                                 |
| <code>aircraft_code</code>          | <code>char(3)</code>   | Aircraft code, IATA                             |
| <code>duration</code>               | <code>interval</code>  | Flight duration                                 |
| <code>days_of_week</code>           | <code>integer[]</code> | Days of the week on which flights are performed |

#### K.4.14. Function `bookings.now`

The demo database contains “snapshots” of data — similar to a backup copy of a real system captured at some point in time. For example, if a flight has the `Departed` status, it means that the aircraft had already departed and was airborne at the time of the backup copy.

The “snapshot” time is saved in the `bookings.now()` function. You can use this function in demo queries for cases where you would use the `now()` function in a real database.

In addition, the return value of this function determines the version of the demo database. The latest version available is of August 15, 2017.

#### K.4.15. Function `bookings.lang`

Some fields in the demo database are available in English and Russian. Translations to other languages are not provided, but are easy to add. The `bookings.lang` returns the value of the `bookings.lang` parameter, that is, the language in which these fields will be displayed.

This function is used in the `aircrafts` and `airports` views and is not intended to be used directly in queries.

### K.5. Usage

#### K.5.1. Schema `bookings`

The `bookings` schema contains all objects of the demo database. When you connect to the database, `search_path` configuration parameter is automatically set to `bookings, public`, so you do not need to specify the schema name explicitly.

However, for the `bookings.now` function, you always have to specify the schema to distinguish this function from the standard `now` function.

#### K.5.2. Translations

By default, values of several translatable fields are shown in Russian. These are `airport_name` and `city` of the `airports` view, as well as `model` of the `aircrafts` view.

You can choose to display these fields in another language (although only the English translation is provided in the demo database). To switch to English, set the `bookings.lang` parameter to `en`. It may be convenient to choose the language at the database level:

```
ALTER DATABASE demo SET bookings.lang = en;
```

You have to reconnect to the database for this command to take effect. For other methods of settings configuration parameters, see [Section 19.1](#).

In the examples below, the English language is selected for translatable fields.

### K.5.3. Sample Queries

To better understand the contents of the demo database, let's take a look at the results of several simple queries.

The results displayed below were received on a small database version (demo-small) of August 15, 2017. If the same queries return different data on your system, check your demo database version (using the `bookings.now` function). Some minor deviations may be caused by the difference between your local time and Moscow time, or your locale settings.

All flights are operated by several types of aircraft:

```
SELECT * FROM aircrafts;
```

| aircraft_code | model               | range |
|---------------|---------------------|-------|
| 773           | Boeing 777-300      | 11100 |
| 763           | Boeing 767-300      | 7900  |
| SU9           | Sukhoi SuperJet-100 | 3000  |
| 320           | Airbus A320-200     | 5700  |
| 321           | Airbus A321-200     | 5600  |
| 319           | Airbus A319-100     | 6700  |
| 733           | Boeing 737-300      | 4200  |
| CN1           | Cessna 208 Caravan  | 1200  |
| CR2           | Bombardier CRJ-200  | 2700  |

(9 rows)

For each aircraft type, a separate list of seats is supported. For example, in a small Cessna 208 Caravan, one can select the following seats:

```
SELECT a.aircraft_code,
 a.model,
 s.seat_no,
 s.fare_conditions
FROM aircrafts a
 JOIN seats s ON a.aircraft_code = s.aircraft_code
WHERE a.model = 'Cessna 208 Caravan'
ORDER BY s.seat_no;
```

| aircraft_code | model              | seat_no | fare_conditions |
|---------------|--------------------|---------|-----------------|
| CN1           | Cessna 208 Caravan | 1A      | Economy         |
| CN1           | Cessna 208 Caravan | 1B      | Economy         |
| CN1           | Cessna 208 Caravan | 2A      | Economy         |
| CN1           | Cessna 208 Caravan | 2B      | Economy         |
| CN1           | Cessna 208 Caravan | 3A      | Economy         |
| CN1           | Cessna 208 Caravan | 3B      | Economy         |
| CN1           | Cessna 208 Caravan | 4A      | Economy         |
| CN1           | Cessna 208 Caravan | 4B      | Economy         |
| CN1           | Cessna 208 Caravan | 5A      | Economy         |
| CN1           | Cessna 208 Caravan | 5B      | Economy         |
| CN1           | Cessna 208 Caravan | 6A      | Economy         |
| CN1           | Cessna 208 Caravan | 6B      | Economy         |

(12 rows)

**Bigger aircraft have more seats of various travel classes:**

```
SELECT s2.aircraft_code,
 string_agg (s2.fare_conditions || '(' || s2.num::text || ')',
 ', ') as fare_conditions
FROM (
 SELECT s.aircraft_code, s.fare_conditions, count(*) as num
 FROM seats s
 GROUP BY s.aircraft_code, s.fare_conditions
 ORDER BY s.aircraft_code, s.fare_conditions
) s2
GROUP BY s2.aircraft_code
ORDER BY s2.aircraft_code;
```

| aircraft_code | fare_conditions                         |
|---------------|-----------------------------------------|
| 319           | Business(20), Economy(96)               |
| 320           | Business(20), Economy(120)              |
| 321           | Business(28), Economy(142)              |
| 733           | Business(12), Economy(118)              |
| 763           | Business(30), Economy(192)              |
| 773           | Business(30), Comfort(48), Economy(324) |
| CN1           | Economy(12)                             |
| CR2           | Economy(50)                             |
| SU9           | Business(12), Economy(85)               |

(9 rows)

**The demo database contains the list of airports of almost all major Russian cities. Most cities have only one airport. The exceptions are:**

```
SELECT a.airport_code as code,
 a.airport_name,
 a.city,
 a.coordinates
FROM airports a
WHERE a.city IN (
 SELECT aa.city
 FROM airports aa
 GROUP BY aa.city
 HAVING COUNT(*) > 1
)
ORDER BY a.city, a.airport_code;
```

| code | airport_name           | city      | coordinates                          |
|------|------------------------|-----------|--------------------------------------|
| DME  | Domodedovo             | Moscow    | (37.9062995910645, 55.4087982177734) |
|      | International Airport  |           |                                      |
| SVO  | Sheremetyevo           | Moscow    | (37.4146, 55.972599)                 |
|      | International Airport  |           |                                      |
| VKO  | Vnukovo                | Moscow    | (37.2615013123, 55.5914993286)       |
|      | International Airport  |           |                                      |
| ULV  | Ulyanovsk              | Ulyanovsk | (48.2266998291, 54.2682991028)       |
|      | Baratayevka Airport    |           |                                      |
| ULY  | Ulyanovsk East Airport | Ulyanovsk | (48.8027000427246, 54.4010009765625) |

(5 rows)

To learn about your flying options from one point to another, it is convenient to use the `routes` materialized view that aggregates information on all flights. For example, here are the destinations where you can get from Volgograd on specific days of the week, with flight duration:

```
SELECT r.arrival_city as city,
 r.arrival_airport as code,
 r.arrival_airport_name as airport_name,
 r.days_of_week,
 r.duration
FROM routes r
WHERE r.departure_city = 'Volgograd';
```

| city        | code | airport_name                       | days_of_week    | duration |
|-------------|------|------------------------------------|-----------------|----------|
| Moscow      | SVO  | Sheremetyevo International Airport | {1,2,3,4,5,6,7} | 01:15:00 |
| Chelyabinsk | CEK  | Chelyabinsk Balandino Airport      | {1,2,3,4,5,6,7} | 01:50:00 |
| Rostov      | ROV  | Rostov-on-Don Airport              | {1,2,3,4,5,6,7} | 00:30:00 |
| Moscow      | VKO  | Vnukovo International Airport      | {1,2,3,4,5,6,7} | 01:10:00 |
| Cheboksary  | CSY  | Cheboksary Airport                 | {1,2,3,4,5,6,7} | 02:45:00 |
| Tomsk       | TOF  | Bogashevo Airport                  | {1}             | 03:50:00 |

(6 rows)

The database was formed at the moment returned by the `bookings.now()` function:

```
SELECT bookings.now() as now;
```

```
now

2017-08-15 18:00:00+03
```

In relation to this moment, all flights are classified as past and future flights:

```
SELECT status,
 count(*) as count,
 min(scheduled_departure) as min_scheduled_departure,
 max(scheduled_departure) as max_scheduled_departure
FROM flights
GROUP BY status
ORDER BY min_scheduled_departure;
```

| status    | count | min_scheduled_departure | max_scheduled_departure |
|-----------|-------|-------------------------|-------------------------|
| Arrived   | 16707 | 2017-07-16 01:50:00+03  | 2017-08-15 17:25:00+03  |
| Cancelled | 414   | 2017-07-19 11:35:00+03  | 2017-09-14 20:55:00+03  |
| Departed  | 58    | 2017-08-15 09:55:00+03  | 2017-08-15 17:50:00+03  |
| Delayed   | 41    | 2017-08-15 15:15:00+03  | 2017-08-16 17:25:00+03  |
| On Time   | 518   | 2017-08-15 17:55:00+03  | 2017-08-16 18:00:00+03  |
| Scheduled | 15383 | 2017-08-16 18:05:00+03  | 2017-09-14 20:40:00+03  |

(6 rows)

Let's find the next flight from Yekaterinburg to Moscow. The `flight` table is not very convenient for such queries, as it does not include information on the cities of departure and arrival. That is why we will use the `flights_v` view:

```
\x
SELECT f.*
FROM flights_v f
WHERE f.departure_city = 'Yekaterinburg'
AND f.arrival_city = 'Moscow'
AND f.scheduled_departure > bookings.now()
ORDER BY f.scheduled_departure
LIMIT 1;
```

```
-[RECORD 1]-----+-----
flight_id | 10927
flight_no | PG0226
scheduled_departure | 2017-08-16 08:10:00+03
scheduled_departure_local | 2017-08-16 10:10:00
scheduled_arrival | 2017-08-16 09:55:00+03
scheduled_arrival_local | 2017-08-16 09:55:00
scheduled_duration | 01:45:00
departure_airport | SVX
departure_airport_name | Koltsovo Airport
departure_city | Yekaterinburg
arrival_airport | SVO
arrival_airport_name | Sheremetyevo International Airport
arrival_city | Moscow
status | On Time
aircraft_code | 773
actual_departure |
actual_departure_local |
actual_arrival |
actual_arrival_local |
actual_duration |
```

Note that the `flights_v` view shows both Moscow time and local time at the airports of departure and arrival.

## K.5.4. Bookings

Each booking can include several tickets, one for each passenger. The ticket, in its turn, can include several flight segments. The complete information about the booking is stored in three tables: `bookings`, `tickets`, and `ticket_flights`.

Let's find several most expensive bookings:

```
SELECT *
FROM bookings
ORDER BY total_amount desc
LIMIT 10;
```

| book_ref | book_date              | total_amount |
|----------|------------------------|--------------|
| 3B54BB   | 2017-07-05 17:08:00+03 | 1204500.00   |
| 3AC131   | 2017-07-31 01:06:00+03 | 1087100.00   |
| 65A6EA   | 2017-07-03 06:28:00+03 | 1065600.00   |
| D7E9AA   | 2017-08-08 05:29:00+03 | 1062800.00   |
| EF479E   | 2017-08-02 15:58:00+03 | 1035100.00   |
| 521C53   | 2017-07-08 09:25:00+03 | 985500.00    |
| 514CA6   | 2017-07-27 05:07:00+03 | 955000.00    |
| D70BD9   | 2017-07-05 12:47:00+03 | 947500.00    |
| EC7EDA   | 2017-07-02 16:13:00+03 | 946800.00    |



```
8E4370 | 2017-07-28 02:04:00+03 | 945700.00
(10 rows)
```

Let's take a look at the tickets included into the booking with code 521C53:

```
SELECT ticket_no,
 passenger_id,
 passenger_name
FROM tickets
WHERE book_ref = '521C53';
```

| ticket_no     | passenger_id | passenger_name     |
|---------------|--------------|--------------------|
| 0005432661914 | 8234 547529  | IVAN IVANOV        |
| 0005432661915 | 2034 201228  | ANTONINA KUZNECOVA |

(2 rows)

If we would like to know, which flight segments are included into Antonina Kuznecova's ticket, we can use the following query:

```
SELECT to_char(f.scheduled_departure, 'DD.MM.YYYY') AS when,
 f.departure_city || ' (' || f.departure_airport || ')' AS departure,
 f.arrival_city || ' (' || f.arrival_airport || ')' AS arrival,
 tf.fare_conditions AS class,
 tf.amount
FROM ticket_flights tf
 JOIN flights_v f ON tf.flight_id = f.flight_id
WHERE tf.ticket_no = '0005432661915'
ORDER BY f.scheduled_departure;
```

| when       | departure             | arrival               | class    | amount    |
|------------|-----------------------|-----------------------|----------|-----------|
| 29.07.2017 | Moscow (SVO)          | Anadyr (DYR)          | Business | 185300.00 |
| 02.08.2017 | Anadyr (DYR)          | Khabarovsk (KHV)      | Business | 92200.00  |
| 03.08.2017 | Khabarovsk (KHV)      | Blagoveshchensk (BQS) | Business | 18000.00  |
| 08.08.2017 | Blagoveshchensk (BQS) | Khabarovsk (KHV)      | Business | 18000.00  |
| 12.08.2017 | Khabarovsk (KHV)      | Anadyr (DYR)          | Economy  | 30700.00  |
| 17.08.2017 | Anadyr (DYR)          | Moscow (SVO)          | Business | 185300.00 |

(6 rows)

As we can see, high booking cost is explained by multiple long-haul flights in business class.

Some of the flight segments in this ticket have earlier dates than the `bookings.now()` return value: it means that these flights had already happened. The last flight had not happened yet at the time of the database creation. After the check-in, a boarding pass with the allocated seat number is issued. We can check the exact seats occupied by Antonina (note the outer left join with table `boarding_passes`):

```
SELECT to_char(f.scheduled_departure, 'DD.MM.YYYY') AS when,
 f.departure_city || ' (' || f.departure_airport || ')' AS departure,
 f.arrival_city || ' (' || f.arrival_airport || ')' AS arrival,
 f.status,
 bp.seat_no
FROM ticket_flights tf
 JOIN flights_v f ON tf.flight_id = f.flight_id
 LEFT JOIN boarding_passes bp ON tf.flight_id = bp.flight_id
 AND tf.ticket_no = bp.ticket_no
WHERE tf.ticket_no = '0005432661915'
```

```
ORDER BY f.scheduled_departure;
```

| when       | departure             | arrival               | status    | seat_no |
|------------|-----------------------|-----------------------|-----------|---------|
| 29.07.2017 | Moscow (SVO)          | Anadyr (DYR)          | Arrived   | 5C      |
| 02.08.2017 | Anadyr (DYR)          | Khabarovsk (KHV)      | Arrived   | 1D      |
| 03.08.2017 | Khabarovsk (KHV)      | Blagoveshchensk (BQS) | Arrived   | 2C      |
| 08.08.2017 | Blagoveshchensk (BQS) | Khabarovsk (KHV)      | Arrived   | 2D      |
| 12.08.2017 | Khabarovsk (KHV)      | Anadyr (DYR)          | Arrived   | 20B     |
| 17.08.2017 | Anadyr (DYR)          | Moscow (SVO)          | Scheduled |         |

(6 rows)

### K.5.5. New Booking

Let's try to send Aleksandr Radishchev from Saint Petersburg to Moscow — the route that made him famous. Naturally, he will travel for free and in business class. We have already found a flight for tomorrow, and a return flight a week later.

```
BEGIN;
```

```
INSERT INTO bookings (book_ref, book_date, total_amount)
VALUES ('_QWE12', bookings.now(), 0);
```

```
INSERT INTO tickets (ticket_no, book_ref, passenger_id, passenger_name)
VALUES ('_0000000000001', '_QWE12', '1749 051790', 'ALEKSANDR RADISHCHEV');
```

```
INSERT INTO ticket_flights (ticket_no, flight_id, fare_conditions, amount)
VALUES ('_0000000000001', 8525, 'Business', 0),
 ('_0000000000001', 4967, 'Business', 0);
```

```
COMMIT;
```

To avoid conflicts with the range of values present in the database, identifiers are started with an underscore.

We will check in Aleksandr for tomorrow's flight right away:

```
INSERT INTO boarding_passes (ticket_no, flight_id, boarding_no, seat_no)
VALUES ('_0000000000001', 8525, 1, '1A');
```

Now let's check the booking information:

```
SELECT b.book_ref,
 t.ticket_no,
 t.passenger_id,
 t.passenger_name,
 tf.fare_conditions,
 tf.amount,
 f.scheduled_departure_local,
 f.scheduled_arrival_local,
 f.departure_city || ' (' || f.departure_airport || ')' AS departure,
 f.arrival_city || ' (' || f.arrival_airport || ')' AS arrival,
 f.status,
 bp.seat_no
FROM bookings b
 JOIN tickets t ON b.book_ref = t.book_ref
 JOIN ticket_flights tf ON tf.ticket_no = t.ticket_no
```

```

 JOIN flights_v f ON tf.flight_id = f.flight_id
 LEFT JOIN boarding_passes bp ON tf.flight_id = bp.flight_id
 AND tf.ticket_no = bp.ticket_no
WHERE b.book_ref = '_QWE12'
ORDER BY t.ticket_no, f.scheduled_departure;

```

```

-[RECORD 1]-----+-----
book_ref | _QWE12
ticket_no | _0000000000001
passenger_id | 1749 051790
passenger_name | ALEKSANDR RADISHCHEV
fare_conditions | Business
amount | 0.00
scheduled_departure_local | 2017-08-16 09:45:00
scheduled_arrival_local | 2017-08-16 10:35:00
departure | St. Petersburg (LED)
arrival | Moscow (SVO)
status | On Time
seat_no | 1A
-[RECORD 2]-----+-----
book_ref | _QWE12
ticket_no | _0000000000001
passenger_id | 1749 051790
passenger_name | ALEKSANDR RADISHCHEV
fare_conditions | Business
amount | 0.00
scheduled_departure_local | 2017-08-23 10:20:00
scheduled_arrival_local | 2017-08-23 11:10:00
departure | Moscow (SVO)
arrival | St. Petersburg (LED)
status | Scheduled
seat_no |

```

We hope that these simple examples helped you get an idea of this demo database.

---

# Appendix L. Acronyms

This is a list of acronyms commonly used in the Postgres Pro documentation and in discussions about Postgres Pro.

AM

*Access Method*

ANSI

*American National Standards Institute*

API

*Application Programming Interface*

ASCII

*American Standard Code for Information Interchange*

BKI

*Backend Interface*

CA

*Certificate Authority*

CIDR

*Classless Inter-Domain Routing*

CPAN

*Comprehensive Perl Archive Network*

CRL

*Certificate Revocation List*

CSV

*Comma Separated Values*

CTE

*Common Table Expression*

CVE

*Common Vulnerabilities and Exposures*

DBA

*Database Administrator*

DBI

*Database Interface (Perl)*

DBMS

*Database Management System*

DDL

*Data Definition Language*, SQL commands such as `CREATE TABLE`, `ALTER USER`

DML

*Data Manipulation Language*, SQL commands such as `INSERT`, `UPDATE`, `DELETE`

DST

*Daylight Saving Time*

ECPG

*Embedded C for Postgres Pro*

ESQL

*Embedded SQL*

FAQ

*Frequently Asked Questions*

FSM

*Free Space Map*

GEQO

*Genetic Query Optimizer*

GIN

*Generalized Inverted Index*

GiST

*Generalized Search Tree*

Git

*Git*

GMT

*Greenwich Mean Time*

GSSAPI

*Generic Security Services Application Programming Interface*

GUC

*Grand Unified Configuration*, the Postgres Pro subsystem that handles server configuration

HBA

*Host-Based Authentication*

HOT

*Heap-Only Tuples*

IEC

*International Electrotechnical Commission*

IEEE

*Institute of Electrical and Electronics Engineers*

IPC

*Inter-Process Communication*

ISO

*International Organization for Standardization*

ISSN

*International Standard Serial Number*

JDBC

*Java Database Connectivity*

JIT

*Just-in-Time compilation*

JSON

*JavaScript Object Notation*

LDAP

*Lightweight Directory Access Protocol*

LSN

*Log Sequence Number*

MCF

Most Common Frequency, that is the frequency associated with some Most Common Value

MCV

Most Common Value, one of the values appearing most often within a particular table column

MITM

*Man-in-the-middle attack*

MSVC

*Microsoft Visual C*

MVCC

*Multi-Version Concurrency Control*

NLS

*National Language Support*

ODBC

*Open Database Connectivity*

OID

*Object Identifier*

OLAP

*Online Analytical Processing*

OLTP

*Online Transaction Processing*

ORDBMS

*Object-Relational Database Management System*

PAM

*Pluggable Authentication Modules*

PGSQL

Postgres Pro

PGXS

Postgres Pro Extension System

PID

*Process Identifier*

PITR

*Point-In-Time Recovery (Continuous Archiving)*

PL

*Procedural Languages (server-side)*

POSIX

*Portable Operating System Interface*

RDBMS

*Relational Database Management System*

RFC

*Request For Comments*

SGML

*Standard Generalized Markup Language*

SNI

*Server Name Indication, RFC 6066*

SPI

*Server Programming Interface*

SP-GiST

*Space-Partitioned Generalized Search Tree*

SQL

*Structured Query Language*

SRF

*Set-Returning Function*

SSH

*Secure Shell*

SSL

*Secure Sockets Layer*

SSPI

*Security Support Provider Interface*

SYSV

*Unix System V*

TCP/IP

*Transmission Control Protocol (TCP) / Internet Protocol (IP)*

TID

*Tuple Identifier*

TLS

*Transport Layer Security*

TOAST

*The Oversized-Attribute Storage Technique*

TPC

*Transaction Processing Performance Council*

URL

*Uniform Resource Locator*

UTC

*Coordinated Universal Time*

UTF

*Unicode Transformation Format*

UTF8

*Eight-Bit Unicode Transformation Format*

UUID

*Universally Unique Identifier*

WAL

*Write-Ahead Log*

XID

*Transaction Identifier*

XML

*Extensible Markup Language*



---

# Appendix M. Glossary

This is a list of terms and their meaning in the context of PostgreSQL and relational database systems in general.

|                              |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
|------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| ACID                         | <a href="#">Atomicity</a> , <a href="#">Consistency</a> , <a href="#">Isolation</a> , and <a href="#">Durability</a> . This set of properties of database transactions is intended to guarantee validity in concurrent operation and even in event of errors, power failures, etc.                                                                                                                                                                                                                                                                                     |
| Aggregate function (routine) | <p>A <a href="#">function</a> that combines (<i>aggregates</i>) multiple input values, for example by counting, averaging or adding, yielding a single output value.</p> <p>For more information, see <a href="#">Section 9.21</a>.</p> <p>See Also <a href="#">Window function (routine)</a>.</p>                                                                                                                                                                                                                                                                     |
| Access Method                | <p>Interfaces which Postgres Pro use in order to access data in tables and indexes. This abstraction allows for adding support for new types of data storage.</p> <p>For more information, see <a href="#">Chapter 62</a> and <a href="#">Chapter 63</a>.</p>                                                                                                                                                                                                                                                                                                          |
| Analytic function            | See <a href="#">Window function (routine)</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
| Analyze (operation)          | <p>The act of collecting statistics from data in <a href="#">tables</a> and other <a href="#">relations</a> to help the <a href="#">query planner</a> to make decisions about how to execute <a href="#">queries</a>.</p> <p>(Don't confuse this term with the <code>ANALYZE</code> option to the <a href="#">EXPLAIN</a> command.)</p> <p>For more information, see <a href="#">ANALYZE</a>.</p>                                                                                                                                                                      |
| Atomic                       | <p>In reference to a <a href="#">datum</a>: the fact that its value cannot be broken down into smaller components.</p> <p>In reference to a <a href="#">database transaction</a>: see <a href="#">atomicity</a>.</p>                                                                                                                                                                                                                                                                                                                                                   |
| Atomicity                    | The property of a <a href="#">transaction</a> that either all its operations complete as a single unit or none do. In addition, if a system failure occurs during the execution of a transaction, no partial results are visible after recovery. This is one of the ACID properties.                                                                                                                                                                                                                                                                                   |
| Attribute                    | An element with a certain name and data type found within a <a href="#">tuple</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Autovacuum (process)         | <p>A set of background processes that routinely perform <a href="#">vacuum</a> and <a href="#">analyze</a> operations. The <a href="#">auxiliary process</a> that coordinates the work and is always present (unless autovacuum is disabled) is known as the <i>autovacuum launcher</i>, and the processes that carry out the tasks are known as the <i>autovacuum workers</i>.</p> <p>For more information, see <a href="#">Section 24.1.6</a>.</p>                                                                                                                   |
| Auxiliary process            | A process within an <a href="#">instance</a> that is in charge of some specific background task for the instance. The auxiliary processes consist of the <a href="#">autovacuum launcher</a> (but not the autovacuum workers), the <a href="#">background writer</a> , the <a href="#">checkpointer</a> , the <a href="#">logger</a> , the <a href="#">startup process</a> , the <a href="#">WAL archiver</a> , the <a href="#">WAL receiver</a> (but not the <a href="#">WAL senders</a> ), the <a href="#">WAL summarizer</a> , and the <a href="#">WAL writer</a> . |
| Backend (process)            | Process of an <a href="#">instance</a> which acts on behalf of a <a href="#">client session</a> and handles its requests.                                                                                                                                                                                                                                                                                                                                                                                                                                              |

|                        |        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|------------------------|--------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                        |        | (Don't confuse this term with the similar terms <i>Background Worker</i> or <i>Background Writer</i> ).                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| Background (process)   | worker | <p>Process within an <i>instance</i>, which runs system- or user-supplied code. Serves as infrastructure for several features in PostgreSQL, such as <i>logical replication</i> and <i>parallel queries</i>. In addition, <i>Extensions</i> can add custom background worker processes.</p> <p>For more information, see <a href="#">Chapter 49</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Background (process)   | writer | <p>An <i>auxiliary process</i> that writes dirty <i>data pages</i> from <i>shared memory</i> to the file system. It wakes up periodically, but works only for a short period in order to distribute its expensive I/O activity over time to avoid generating larger I/O peaks which could block other processes.</p> <p>For more information, see <a href="#">Section 19.4.5</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| Bloat                  |        | Space in data pages which does not contain current row versions, such as unused (free) space or outdated row versions.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Bootstrap superuser    |        | <p>The first <i>user</i> initialized in a <i>database cluster</i>.</p> <p>This user owns all system catalog tables in each database. It is also the role from which all granted permissions originate. Because of these things, this role may not be dropped.</p> <p>This role also behaves as a normal <i>database superuser</i>, and its superuser status cannot be removed.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| Buffer Access Strategy |        | <p>Some operations will access a large number of <i>pages</i>. A <i>Buffer Access Strategy</i> helps to prevent these operations from evicting too many pages from <i>shared buffers</i>.</p> <p>A Buffer Access Strategy sets up references to a limited number of <i>shared buffers</i> and reuses them circularly. When the operation requires a new page, a victim buffer is chosen from the buffers in the strategy ring, which may require flushing the page's dirty data and possibly also unflushed <i>WAL</i> to permanent storage.</p> <p>Buffer Access Strategies are used for various operations such as sequential scans of large tables, <code>VACUUM</code>, <code>COPY</code>, <code>CREATE TABLE AS SELECT</code>, <code>ALTER TABLE</code>, <code>CREATE DATABASE</code>, <code>CREATE INDEX</code>, and <code>CLUSTER</code>.</p> |
| Cast                   |        | <p>A conversion of a <i>datum</i> from its current data type to another data type.</p> <p>For more information, see <a href="#">CREATE CAST</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Catalog                |        | <p>The SQL standard uses this term to indicate what is called a <i>database</i> in PostgreSQL's terminology.</p> <p>(Don't confuse this term with <i>system catalog</i>).</p> <p>For more information, see <a href="#">Section 22.1</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
| Check constraint       |        | <p>A type of <i>constraint</i> defined on a <i>relation</i> which restricts the values allowed in one or more <i>attributes</i>. The check constraint can make reference to any attribute of the same row in the relation, but cannot reference other rows of the same relation or other relations.</p> <p>For more information, see <a href="#">Section 5.5</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |

|                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
|------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Checkpoint             | <p>A point in the <a href="#">WAL</a> sequence at which it is guaranteed that the heap and index data files have been updated with all information from <a href="#">shared memory</a> modified before that checkpoint; a <i>checkpoint record</i> is written and flushed to WAL to mark that point.</p> <p>A checkpoint is also the act of carrying out all the actions that are necessary to reach a checkpoint as defined above. This process is initiated when pre-defined conditions are met, such as a specified amount of time has passed, or a certain volume of records has been written; or it can be invoked by the user with the command <code>CHECKPOINT</code>.</p> <p>For more information, see <a href="#">Section 28.5</a>.</p> |
| Checkpointer (process) | An <a href="#">auxiliary process</a> that is responsible for executing <i>checkpoints</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Class (archaic)        | See <a href="#">Relation</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Client (process)       | Any process, possibly remote, that establishes a <i>session</i> by <i>connecting</i> to an <i>instance</i> to interact with a <i>database</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Cluster owner          | <p>The operating system user that owns the <a href="#">data directory</a> and under which the <code>postgres</code> process is run. It is required that this user exist prior to creating a new <i>database cluster</i>.</p> <p>On operating systems with a <code>root</code> user, said user is not allowed to be the cluster owner.</p>                                                                                                                                                                                                                                                                                                                                                                                                       |
| Column                 | An <i>attribute</i> found in a <i>table</i> or <i>view</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Commit                 | <p>The act of finalizing a <i>transaction</i> within the <i>database</i>, which makes it visible to other transactions and assures its <i>durability</i>.</p> <p>For more information, see <a href="#">COMMIT</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Concurrency            | The concept that multiple independent operations happen within the <i>database</i> at the same time. In PostgreSQL, concurrency is controlled by the <i>multiversion concurrency control</i> mechanism.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
| Connection             | <p>An established line of communication between a client process and a <i>back-end</i> process, usually over a network, supporting a <i>session</i>. This term is sometimes used as a synonym for session.</p> <p>For more information, see <a href="#">Section 19.3</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| Consistency            | The property that the data in the <i>database</i> is always in compliance with <i>integrity constraints</i> . Transactions may be allowed to violate some of the constraints transiently before it commits, but if such violations are not resolved by the time it commits, such a transaction is automatically <i>rolled back</i> . This is one of the ACID properties.                                                                                                                                                                                                                                                                                                                                                                        |
| Constraint             | <p>A restriction on the values of data allowed within a <i>table</i>, or in attributes of a <i>domain</i>.</p> <p>For more information, see <a href="#">Section 5.5</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Cumulative System      | <p>Statistics</p> <p>A system which, if enabled, accumulates statistical information about the <i>instance</i>'s activities.</p> <p>For more information, see <a href="#">Section 27.2</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Data area              | See <a href="#">Data directory</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |

|                    |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
|--------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Database           | <p>A named collection of <i>local SQL objects</i>.</p> <p>For more information, see <a href="#">Section 22.1</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Database cluster   | <p>A collection of databases and global SQL objects, and their common static and dynamic metadata. Sometimes referred to as a <i>cluster</i>. A database cluster is created using the <code>initdb</code> program.</p> <p>In PostgreSQL, the term <i>cluster</i> is also sometimes used to refer to an instance. (Don't confuse this term with the SQL command <code>CLUSTER</code>.)</p> <p>See also <i>cluster owner</i>, the operating-system owner of a cluster, and <i>bootstrap superuser</i>, the PostgreSQL owner of a cluster.</p> |
| Database server    | See <a href="#">Instance</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| Database superuser | <p>A role having <i>superuser status</i> (see <a href="#">Section 21.2</a>).</p> <p>Frequently referred to as <i>superuser</i>.</p>                                                                                                                                                                                                                                                                                                                                                                                                         |
| Data directory     | <p>The base directory on the file system of a <i>server</i> that contains all data files and subdirectories associated with a <i>database cluster</i> (with the exception of <i>tablespaces</i>, and optionally <i>WAL</i>). The environment variable <code>PGDATA</code> is commonly used to refer to the data directory.</p> <p>A <i>cluster's</i> storage space comprises the data directory plus any additional tablespaces.</p> <p>For more information, see <a href="#">Section 66.1</a>.</p>                                         |
| Data page          | <p>The basic structure used to store relation data. All pages are of the same size. Data pages are typically stored on disk, each in a specific file, and can be read to <i>shared buffers</i> where they can be modified, becoming <i>dirty</i>. They become clean when written to disk. New pages, which initially exist in memory only, are also dirty until written.</p>                                                                                                                                                                |
| Datum              | The internal representation of one value of an SQL data type.                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Delete             | <p>An SQL command which removes <i>rows</i> from a given <i>table</i> or <i>relation</i>.</p> <p>For more information, see <a href="#">DELETE</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                      |
| Domain             | <p>A user-defined data type that is based on another underlying data type. It acts the same as the underlying type except for possibly restricting the set of allowed values.</p> <p>For more information, see <a href="#">Section 8.18</a>.</p>                                                                                                                                                                                                                                                                                            |
| Durability         | The assurance that once a <i>transaction</i> has been <i>committed</i> , the changes remain even after a system failure or crash. This is one of the ACID properties.                                                                                                                                                                                                                                                                                                                                                                       |
| Extension          | <p>A software add-on package that can be installed on an <i>instance</i> to get extra features.</p> <p>For more information, see <a href="#">Section 39.17</a>.</p>                                                                                                                                                                                                                                                                                                                                                                         |
| File segment       | <p>A physical file which stores data for a given <i>relation</i>. File segments are limited in size by a configuration value (typically 1 gigabyte), so if a relation exceeds that size, it is split into multiple segments.</p> <p>For more information, see <a href="#">Section 66.1</a>.</p>                                                                                                                                                                                                                                             |

|                          |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
|--------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                          | (Don't confuse this term with the similar term <a href="#">WAL segment</a> ).                                                                                                                                                                                                                                                                                                                                                                                                       |
| Foreign data wrapper     | A means of representing data that is not contained in the local <a href="#">database</a> so that it appears as if were in local <a href="#">table(s)</a> . With a foreign data wrapper it is possible to define a <a href="#">foreign server</a> and <a href="#">foreign tables</a> .                                                                                                                                                                                               |
| Foreign key              | A type of <a href="#">constraint</a> defined on one or more <a href="#">columns</a> in a <a href="#">table</a> which requires the value(s) in those <a href="#">columns</a> to identify zero or one <a href="#">row</a> in another (or, infrequently, the same) <a href="#">table</a> .                                                                                                                                                                                             |
| Foreign server           | <p>A named collection of <a href="#">foreign tables</a> which all use the same <a href="#">foreign data wrapper</a> and have other configuration values in common.</p> <p>For more information, see <a href="#">CREATE SERVER</a>.</p>                                                                                                                                                                                                                                              |
| Foreign table (relation) | A <a href="#">relation</a> which appears to have <a href="#">rows</a> and <a href="#">columns</a> similar to a regular <a href="#">table</a> , but will forward requests for data through its <a href="#">foreign data wrapper</a> , which will return <a href="#">result sets</a> structured according to the definition of the <a href="#">foreign table</a> .                                                                                                                    |
| Fork                     | Each of the separate segmented file sets in which a relation is stored. The <i>main fork</i> is where the actual data resides. There also exist two secondary forks for metadata: the <a href="#">free space map</a> and the <a href="#">visibility map</a> . <a href="#">Unlogged relations</a> also have an <i>init fork</i> .                                                                                                                                                    |
| Free space map (fork)    | <p>A storage structure that keeps metadata about each data page of a table's main fork. The free space map entry for each page stores the amount of free space that's available for future tuples, and is structured to be efficiently searched for available space for a new tuple of a given size.</p> <p>For more information, see <a href="#">Section 66.3</a>.</p>                                                                                                             |
| Function (routine)       | <p>A type of routine that receives zero or more arguments, returns zero or more output values, and is constrained to run within one transaction. Functions are invoked as part of a query, for example via <code>SELECT</code>. Certain functions can return <a href="#">sets</a>; those are called <i>set-returning functions</i>.</p> <p>Functions can also be used for <a href="#">triggers</a> to invoke.</p> <p>For more information, see <a href="#">CREATE FUNCTION</a>.</p> |
| GMT                      | See <a href="#">UTC</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Grant                    | <p>An SQL command that is used to allow a <a href="#">user</a> or <a href="#">role</a> to access specific objects within the <a href="#">database</a>.</p> <p>For more information, see <a href="#">GRANT</a>.</p>                                                                                                                                                                                                                                                                  |
| Heap                     | Contains the values of <a href="#">row</a> attributes (i.e., the data) for a <a href="#">relation</a> . The heap is realized within one or more <a href="#">file segments</a> in the relation's <i>main fork</i> .                                                                                                                                                                                                                                                                  |
| Host                     | A computer that communicates with other computers over a network. This is sometimes used as a synonym for <a href="#">server</a> . It is also used to refer to a computer where <a href="#">client processes</a> run.                                                                                                                                                                                                                                                               |
| Index (relation)         | <p>A <a href="#">relation</a> that contains data derived from a <a href="#">table</a> or <a href="#">materialized view</a>. Its internal structure supports fast retrieval of and access to the original data.</p> <p>For more information, see <a href="#">CREATE INDEX</a>.</p>                                                                                                                                                                                                   |
| Insert                   | An SQL command used to add new data into a <a href="#">table</a> .                                                                                                                                                                                                                                                                                                                                                                                                                  |

For more information, see [INSERT](#).

|                           |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
|---------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Instance                  | <p>A group of <a href="#">backend</a> and <a href="#">auxiliary processes</a> that communicate using a common shared memory area. One <a href="#">postmaster process</a> manages the instance; one instance manages exactly one <a href="#">database cluster</a> with all its databases. Many instances can run on the same <a href="#">server</a> as long as their TCP ports do not conflict.</p> <p>The instance handles all key features of a DBMS: read and write access to files and shared memory, assurance of the ACID properties, <a href="#">connections</a> to <a href="#">client processes</a>, privilege verification, crash recovery, replication, etc.</p> |
| Isolation                 | <p>The property that the effects of a transaction are not visible to <a href="#">concurrent transactions</a> before it commits. This is one of the ACID properties.</p> <p>For more information, see <a href="#">Section 13.2</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Join                      | <p>An operation and SQL keyword used in <a href="#">queries</a> for combining data from multiple <a href="#">relations</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Key                       | <p>A means of identifying a <a href="#">row</a> within a <a href="#">table</a> or other <a href="#">relation</a> by values contained within one or more <a href="#">attributes</a> in that relation.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Lock                      | <p>A mechanism that allows a process to limit or prevent simultaneous access to a resource.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Log file                  | <p>Log files contain human-readable text lines about events. Examples include login failures, long-running queries, etc.</p> <p>For more information, see <a href="#">Section 24.4</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Logged                    | <p>A <a href="#">table</a> is considered <a href="#">logged</a> if changes to it are sent to the <a href="#">WAL</a>. By default, all regular tables are logged. A table can be specified as <a href="#">unlogged</a> either at creation time or via the <code>ALTER TABLE</code> command.</p>                                                                                                                                                                                                                                                                                                                                                                            |
| Logger (process)          | <p>An <a href="#">auxiliary process</a> which, if enabled, writes information about database events into the current <a href="#">log file</a>. When reaching certain time- or volume-dependent criteria, a new log file is created. Also called <i>syslogger</i>.</p> <p>For more information, see <a href="#">Section 19.8</a>.</p>                                                                                                                                                                                                                                                                                                                                      |
| Log record                | <p>Archaic term for a <a href="#">WAL record</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Log sequence number (LSN) | <p>Byte offset into the <a href="#">WAL</a>, increasing monotonically with each new <a href="#">WAL record</a>.</p> <p>For more information, see <a href="#">pg_lsn</a> and <a href="#">Section 28.7</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| LSN                       | <p>See <a href="#">Log sequence number</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Master (server)           | <p>See <a href="#">Primary (server)</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| Materialized              | <p>The property that some information has been pre-computed and stored for later use, rather than computing it on-the-fly.</p> <p>This term is used in <a href="#">materialized view</a>, to mean that the data derived from the view's query is stored on disk separately from the sources of that data.</p> <p>This term is also used to refer to some multi-step queries to mean that the data resulting from executing a given step is stored in memory (with the</p>                                                                                                                                                                                                 |

|                                          |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
|------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                          | possibility of spilling to disk), so that it can be read multiple times by another step.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Materialized view (relation)             | <p>A <i>relation</i> that is defined by a <code>SELECT</code> statement (just like a <i>view</i>), but stores data in the same way that a <i>table</i> does. It cannot be modified via <code>INSERT</code>, <code>UPDATE</code>, <code>DELETE</code>, or <code>MERGE</code> operations.</p> <p>For more information, see <a href="#">CREATE MATERIALIZED VIEW</a>.</p>                                                                                                                                                                                                                                                                                       |
| Merge                                    | <p>An SQL command used to conditionally add, modify, or remove <i>rows</i> in a given <i>table</i>, using data from a source <i>relation</i>.</p> <p>For more information, see <a href="#">MERGE</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Multi-version concurrency control (MVCC) | A mechanism designed to allow several <i>transactions</i> to be reading and writing the same rows without one process causing other processes to stall. In PostgreSQL, MVCC is implemented by creating copies ( <i>versions</i> ) of <i>tuples</i> as they are modified; after transactions that can see the old versions terminate, those old versions need to be removed.                                                                                                                                                                                                                                                                                  |
| Null                                     | A concept of non-existence that is a central tenet of relational database theory. It represents the absence of a definite value.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Optimizer                                | See <a href="#">Query planner</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Parallel query                           | The ability to handle parts of executing a <i>query</i> to take advantage of parallel processes on servers with multiple CPUs.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Partition                                | <p>One of several disjoint (not overlapping) subsets of a larger set.</p> <p>In reference to a <i>partitioned table</i>: One of the tables that each contain part of the data of the partitioned table, which is said to be the <i>parent</i>. The partition is itself a table, so it can also be queried directly; at the same time, a partition can sometimes be a partitioned table, allowing hierarchies to be created.</p> <p>In reference to a <i>window function</i> in a <i>query</i>, a partition is a user-defined criterion that identifies which neighboring <i>rows</i> of the <i>query's result set</i> can be considered by the function.</p> |
| Partitioned table (relation)             | A <i>relation</i> that is in semantic terms the same as a <i>table</i> , but whose storage is distributed across several <i>partitions</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Postmaster (process)                     | <p>The very first process of an <i>instance</i>. It starts and manages the <i>auxiliary processes</i> and creates <i>backend processes</i> on demand.</p> <p>For more information, see <a href="#">Section 18.3</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Primary key                              | A special case of a <i>unique constraint</i> defined on a <i>table</i> or other <i>relation</i> that also guarantees that all of the <i>attributes</i> within the <i>primary key</i> do not have <i>null</i> values. As the name implies, there can be only one primary key per table, though it is possible to have multiple unique constraints that also have no null-capable attributes.                                                                                                                                                                                                                                                                  |
| Primary (server)                         | When two or more <i>databases</i> are linked via <i>replication</i> , the <i>server</i> that is considered the authoritative source of information is called the <i>primary</i> , also known as a <i>master</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| Procedure (routine)                      | A type of routine. Their distinctive qualities are that they do not return values, and that they are allowed to make transactional statements such as <code>COMMIT</code> and <code>ROLLBACK</code> . They are invoked via the <code>CALL</code> command.                                                                                                                                                                                                                                                                                                                                                                                                    |



|                       |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
|-----------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                       | For more information, see <a href="#">CREATE PROCEDURE</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| Query                 | A request sent by a client to a <a href="#">backend</a> , usually to return results or to modify data on the database.                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Query planner         | The part of PostgreSQL that is devoted to determining ( <i>planning</i> ) the most efficient way to execute <a href="#">queries</a> . Also known as <i>query optimizer</i> , <i>optimizer</i> , or simply <i>planner</i> .                                                                                                                                                                                                                                                                                                                                              |
| Record                | See <a href="#">Tuple</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Recycling             | See <a href="#">WAL file</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Referee (node)        | A node that only participates in voting and does not contain any user data.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Referential integrity | A means of restricting data in one <a href="#">relation</a> by a <a href="#">foreign key</a> so that it must have matching data in another <a href="#">relation</a> .                                                                                                                                                                                                                                                                                                                                                                                                   |
| Relation              | <p>The generic term for all objects in a <a href="#">database</a> that have a name and a list of <a href="#">attributes</a> defined in a specific order. <a href="#">Tables</a>, <a href="#">sequences</a>, <a href="#">views</a>, <a href="#">foreign tables</a>, <a href="#">materialized views</a>, composite types, and <a href="#">indexes</a> are all relations.</p> <p>More generically, a relation is a set of tuples; for example, the result of a query is also a relation.</p> <p>In PostgreSQL, <i>Class</i> is an archaic synonym for <i>relation</i>.</p> |
| Replica (server)      | A <a href="#">database</a> that is paired with a <a href="#">primary</a> database and is maintaining a copy of some or all of the primary database's data. The foremost reasons for doing this are to allow for greater access to that data, and to maintain availability of the data in the event that the <a href="#">primary</a> becomes unavailable.                                                                                                                                                                                                                |
| Replication           | The act of reproducing data on one <a href="#">server</a> onto another server called a <a href="#">replica</a> . This can take the form of <i>physical replication</i> , where all file changes from one server are copied verbatim, or <i>logical replication</i> where a defined subset of data changes are conveyed using a higher-level representation.                                                                                                                                                                                                             |
| Restartpoint          | <p>A variant of a <a href="#">checkpoint</a> performed on a <a href="#">replica</a>.</p> <p>For more information, see <a href="#">Section 28.5</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Result set            | <p>A <a href="#">relation</a> transmitted from a <a href="#">backend process</a> to a <a href="#">client</a> upon the completion of an SQL command, usually a <code>SELECT</code> but it can be an <code>INSERT</code>, <code>UPDATE</code>, <code>DELETE</code>, or <code>MERGE</code> command if the <code>RETURNING</code> clause is specified.</p> <p>The fact that a result set is a relation means that a query can be used in the definition of another query, becoming a <i>subquery</i>.</p>                                                                   |
| Revoke                | <p>A command to prevent access to a named set of <a href="#">database</a> objects for a named list of <a href="#">roles</a>.</p> <p>For more information, see <a href="#">REVOKE</a>.</p>                                                                                                                                                                                                                                                                                                                                                                               |
| Role                  | <p>A collection of access privileges to the <a href="#">instance</a>. Roles are themselves a privilege that can be granted to other roles. This is often done for convenience or to ensure completeness when multiple <a href="#">users</a> need the same privileges.</p> <p>For more information, see <a href="#">CREATE ROLE</a>.</p>                                                                                                                                                                                                                                 |



|                     |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
|---------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Rollback            | <p>A command to undo all of the operations performed since the beginning of a <a href="#">transaction</a>.</p> <p>For more information, see <a href="#">ROLLBACK</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Routine             | <p>A defined set of instructions stored in the database system that can be invoked for execution. A routine can be written in a variety of programming languages. Routines can be <a href="#">functions</a> (including set-returning functions and <a href="#">trigger functions</a>), <a href="#">aggregate functions</a>, and <a href="#">procedures</a>.</p> <p>Many routines are already defined within PostgreSQL itself, but user-defined ones can also be added.</p>                                                                                                                                                               |
| Row                 | <p>See <a href="#">Tuple</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
| Savepoint           | <p>A special mark in the sequence of steps in a <a href="#">transaction</a>. Data modifications after this point in time may be reverted to the time of the savepoint.</p> <p>For more information, see <a href="#">SAVEPOINT</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                    |
| Schema              | <p>A schema is a namespace for <a href="#">SQL objects</a>, which all reside in the same <a href="#">database</a>. Each SQL object must reside in exactly one schema.</p> <p>All system-defined SQL objects reside in schema <code>pg_catalog</code>.</p> <p>More generically, the term <i>schema</i> is used to mean all data descriptions (<a href="#">table</a> definitions, <a href="#">constraints</a>, comments, etc.) for a given <a href="#">database</a> or subset thereof.</p> <p>For more information, see <a href="#">Section 5.10</a>.</p>                                                                                   |
| Segment             | <p>See <a href="#">File segment</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Select              | <p>The SQL command used to request data from a <a href="#">database</a>. Normally, <code>SELECT</code> commands are not expected to modify the <a href="#">database</a> in any way, but it is possible that <a href="#">functions</a> invoked within the query could have side effects that do modify data.</p> <p>For more information, see <a href="#">SELECT</a>.</p>                                                                                                                                                                                                                                                                  |
| Sequence (relation) | <p>A type of relation that is used to generate values. Typically the generated values are sequential non-repeating numbers. They are commonly used to generate surrogate <a href="#">primary key</a> values.</p>                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Server              | <p>A computer on which PostgreSQL <a href="#">instances</a> run. The term <i>server</i> denotes real hardware, a container, or a <i>virtual machine</i>.</p> <p>This term is sometimes used to refer to an instance or to a host.</p>                                                                                                                                                                                                                                                                                                                                                                                                     |
| Session             | <p>A state that allows a client and a backend to interact, communicating over a <a href="#">connection</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Shared memory       | <p>RAM which is used by the processes common to an <a href="#">instance</a>. It mirrors parts of <a href="#">database</a> files, provides a transient area for <a href="#">WAL records</a>, and stores additional common information. Note that shared memory belongs to the complete instance, not to a single database.</p> <p>The largest part of shared memory is known as <i>shared buffers</i> and is used to mirror part of data files, organized into pages. When a page is modified, it is called a dirty page until it is written back to the file system.</p> <p>For more information, see <a href="#">Section 19.4.1</a>.</p> |

|                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
|------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| SQL object       | <p>Any object that can be created with a <code>CREATE</code> command. Most objects are specific to one database, and are commonly known as <i>local objects</i>.</p> <p>Most local objects reside in a specific <i>schema</i> in their containing database, such as <i>relations</i> (all types), <i>routines</i> (all types), data types, etc. The names of such objects of the same type in the same schema are enforced to be unique.</p> <p>There also exist local objects that do not reside in schemas; some examples are <i>extensions</i>, <i>data type casts</i>, and <i>foreign data wrappers</i>. The names of such objects of the same type are enforced to be unique within the database.</p> <p>Other object types, such as <i>roles</i>, <i>tablespaces</i>, replication origins, subscriptions for logical replication, and databases themselves are not local SQL objects since they exist entirely outside of any specific database; they are called <i>global objects</i>. The names of such objects are enforced to be unique within the whole database cluster.</p> <p>For more information, see <a href="#">Section 22.1</a>.</p> |
| SQL standard     | A series of documents that define the SQL language.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Standby (server) | See <a href="#">Replica (server)</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Startup process  | <p>An <i>auxiliary process</i> that replays WAL during crash recovery and in a <i>physical replica</i>.</p> <p>(The name is historical: the startup process was named before replication was implemented; the name refers to its task as it relates to the server startup following a crash.)</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Superuser        | As used in this documentation, it is a synonym for <i>database superuser</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| System catalog   | <p>A collection of <i>tables</i> which describe the structure of all <i>SQL objects</i> of the instance. The system catalog resides in the schema <code>pg_catalog</code>. These tables contain data in internal representation and are not typically considered useful for user examination; a number of user-friendlier <i>views</i>, also in schema <code>pg_catalog</code>, offer more convenient access to some of that information, while additional tables and views exist in schema <code>information_schema</code> (see <a href="#">Chapter 38</a>) that expose some of the same and additional information as mandated by the <i>SQL standard</i>.</p> <p>For more information, see <a href="#">Section 5.10</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                         |
| Table            | <p>A collection of <i>tuples</i> having a common data structure (the same number of <i>attributes</i>, in the same order, having the same name and type per position). A table is the most common form of <i>relation</i> in PostgreSQL.</p> <p>For more information, see <a href="#">CREATE TABLE</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Tablespace       | <p>A named location on the server file system. All <i>SQL objects</i> which require storage beyond their definition in the <i>system catalog</i> must belong to a single tablespace. Initially, a database cluster contains a single usable tablespace which is used as the default for all SQL objects, called <code>pg_default</code>.</p> <p>For more information, see <a href="#">Section 22.7</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Temporary table  | <i>Tables</i> that exist either for the lifetime of a <i>session</i> or a <i>transaction</i> , as specified at the time of creation. The data in them is not visible to other sessions, and is not <i>logged</i> . Temporary tables are often used to store intermediate data for a multi-step operation.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |

|                               |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
|-------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                               | For more information, see <a href="#">CREATE TABLE</a> .                                                                                                                                                                                                                                                                                                                                                                                                                       |
| TOAST                         | <p>A mechanism by which large attributes of table rows are split and stored in a secondary table, called the <i>TOAST table</i>. Each relation with large attributes has its own TOAST table.</p> <p>For more information, see <a href="#">Section 66.2</a>.</p>                                                                                                                                                                                                               |
| Transaction                   | <p>A combination of commands that must act as a single <i>atomic</i> command: they all succeed or all fail as a single unit, and their effects are not visible to other <i>sessions</i> until the transaction is complete, and possibly even later, depending on the isolation level.</p> <p>For more information, see <a href="#">Section 13.2</a>.</p>                                                                                                                       |
| Transaction ID                | <p>The numerical, unique, sequentially-assigned identifier that each transaction receives when it first causes a database modification. Frequently abbreviated as <i>xid</i>.</p> <p>For more information, see <a href="#">Section 8.19</a>.</p>                                                                                                                                                                                                                               |
| Transactions per second (TPS) | <p>Average number of transactions that are executed per second, totaled across all sessions active for a measured run. This is used as a measure of the performance characteristics of an instance.</p>                                                                                                                                                                                                                                                                        |
| Trigger                       | <p>A <i>function</i> which can be defined to execute whenever a certain operation (INSERT, UPDATE, DELETE, TRUNCATE) is applied to a <i>relation</i>. A trigger executes within the same <i>transaction</i> as the statement which invoked it, and if the function fails, then the invoking statement also fails.</p> <p>For more information, see <a href="#">CREATE TRIGGER</a>.</p>                                                                                         |
| Tuple                         | <p>A collection of <i>attributes</i> in a fixed order. That order may be defined by the <i>table</i> (or other <i>relation</i>) where the tuple is contained, in which case the tuple is often called a <i>row</i>. It may also be defined by the structure of a result set, in which case it is sometimes called a <i>record</i>.</p>                                                                                                                                         |
| Unique constraint             | <p>A type of <i>constraint</i> defined on a <i>relation</i> which restricts the values allowed in one or a combination of columns so that each value or combination of values can only appear once in the relation — that is, no other row in the relation contains values that are equal to those.</p> <p>Because <i>null values</i> are not considered equal to each other, multiple rows with null values are allowed to exist without violating the unique constraint.</p> |
| Unlogged                      | <p>The property of certain <i>relations</i> that the changes to them are not reflected in the <i>WAL</i>. This disables replication and crash recovery for these relations.</p> <p>The primary use of unlogged tables is for storing transient work data that must be shared across processes.</p> <p><i>Temporary tables</i> are always unlogged.</p>                                                                                                                         |
| Update                        | <p>An SQL command used to modify <i>rows</i> that may already exist in a specified <i>table</i>. It cannot create or remove rows.</p> <p>For more information, see <a href="#">UPDATE</a>.</p>                                                                                                                                                                                                                                                                                 |
| User                          | <p>A <i>role</i> that has the <i>login privilege</i> (see <a href="#">Section 21.2</a>).</p>                                                                                                                                                                                                                                                                                                                                                                                   |

|                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| User mapping           | The translation of login credentials in the local <i>database</i> to credentials in a remote data system defined by a <i>foreign data wrapper</i> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
| UTC                    | Universal Coordinated Time, the primary global time reference, approximately the time prevailing at the zero meridian of longitude. Often but inaccurately referred to as GMT (Greenwich Mean Time).                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Vacuum                 | <p>The process of removing outdated <i>tuple versions</i> from tables or materialized views, and other closely related processing required by PostgreSQL's implementation of <i>MVCC</i>. This can be initiated through the use of the <code>VACUUM</code> command, but can also be handled automatically via <i>autovacuum</i> processes.</p> <p>For more information, see <a href="#">Section 24.1</a> .</p>                                                                                                                                                                                                                                                                                                                                                                                                             |
| View                   | <p>A <i>relation</i> that is defined by a <code>SELECT</code> statement, but has no storage of its own. Any time a query references a view, the definition of the view is substituted into the query as if the user had typed it as a subquery instead of the name of the view.</p> <p>For more information, see <a href="#">CREATE VIEW</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Visibility map (fork)  | A storage structure that keeps metadata about each data page of a table's main fork. The visibility map entry for each page stores two bits: the first one ( <i>all-visible</i> ) indicates that all tuples in the page are visible to all transactions. The second one ( <i>all-frozen</i> ) indicates that all tuples in the page are marked frozen.                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| WAL                    | See <a href="#">Write-ahead log</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| WAL archiver (process) | <p>An <i>auxiliary process</i> which, if enabled, saves copies of <i>WAL files</i> for the purpose of creating backups or keeping <i>replicas</i> current.</p> <p>For more information, see <a href="#">Section 25.1</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| WAL file               | <p>Also known as <i>WAL segment</i> or <i>WAL segment file</i>. Each of the sequentially-numbered files that provide storage space for <i>WAL</i>. The files are all of the same predefined size and are written in sequential order, interspersing changes as they occur in multiple simultaneous sessions. If the system crashes, the files are read in order, and each of the changes is replayed to restore the system to the state it was in before the crash.</p> <p>Each WAL file can be released after a <i>checkpoint</i> writes all the changes in it to the corresponding data files. Releasing the file can be done either by deleting it, or by changing its name so that it will be used in the future, which is called <i>recycling</i>.</p> <p>For more information, see <a href="#">Section 28.7</a>.</p> |
| WAL record             | <p>A low-level description of an individual data change. It contains sufficient information for the data change to be re-executed (<i>replayed</i>) in case a system failure causes the change to be lost. WAL records use a non-printable binary format.</p> <p>For more information, see <a href="#">Section 28.7</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| WAL receiver (process) | <p>An <i>auxiliary process</i> that runs on a <i>replica</i> to receive WAL from the <i>primary server</i> for replay by the <i>startup process</i>.</p> <p>For more information, see <a href="#">Section 26.2</a>.</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| WAL segment            | See <a href="#">WAL file</a> .                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |

|                           |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
|---------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| WAL sender (process)      | A special <i>backend process</i> that streams WAL over a network. The receiving end can be a <i>WAL receiver</i> in a <i>replica</i> , <i>pg_receivewal</i> , or any other client program that speaks the replication protocol.                                                                                                                                                                                                                                                                              |
| WAL summarizer (process)  | An <i>auxiliary process</i> that summarizes WAL data .<br><br>For more information, see <a href="#">Section 19.5.7</a> .                                                                                                                                                                                                                                                                                                                                                                                     |
| WAL writer (process)      | An <i>auxiliary process</i> that writes <i>WAL records</i> from <i>shared memory</i> to <i>WAL files</i> .<br><br>For more information, see <a href="#">Section 19.5</a> .                                                                                                                                                                                                                                                                                                                                   |
| Window function (routine) | A type of <i>function</i> used in a <i>query</i> that applies to a <i>partition</i> of the query's <i>result set</i> ; the function's result is based on values found in <i>rows</i> of the same partition or frame.<br><br>All <i>aggregate functions</i> can be used as window functions, but window functions can also be used to, for example, give ranks to each of the rows in the partition. Also known as <i>analytic functions</i> .<br><br>For more information, see <a href="#">Section 3.5</a> . |
| Write-ahead log           | The journal that keeps track of the changes in the <i>database cluster</i> as user- and system-invoked operations take place. It comprises many individual <i>WAL records</i> written sequentially to <i>WAL files</i> .                                                                                                                                                                                                                                                                                     |

---

# Appendix N. Color Support

Most programs in the Postgres Pro package can produce colorized console output. This appendix describes how that is configured.

## N.1. When Color is Used

To use colorized output, set the environment variable `PG_COLORS` as follows:

1. If the value is `always`, then color is used.
2. If the value is `auto` and the standard error stream is associated with a terminal device, then color is used.
3. Otherwise, color is not used.

## N.2. Configuring the Colors

The actual colors to be used are configured using the environment variable `PG_COLORS` (note plural). The value is a colon-separated list of *key=value* pairs. The keys specify what the color is to be used for. The values are SGR (Select Graphic Rendition) specifications, which are interpreted by the terminal.

The following keys are currently in use:

`error`

used to highlight the text “error” in error messages

`warning`

used to highlight the text “warning” in warning messages

`note`

used to highlight the text “detail” and “hint” in such messages

`locus`

used to highlight location information (e.g., program name and file name) in messages

The default value is `error=01;31:warning=01;35:note=01;36:locus=01` (`01;31` = bold red, `01;35` = bold magenta, `01;36` = bold cyan, `01` = bold default color).

### Tip

This color specification format is also used by other software packages such as GCC, GNU coreutils, and GNU grep.

---

# Appendix O. Obsolete or Renamed Features

Functionality is sometimes removed from PostgreSQL, feature, setting and file names sometimes change, or documentation moves to different places. This section directs users coming from old versions of the documentation or from external links to the appropriate new location for the information they need.

## O.1. `recovery.conf` file merged into `postgresql.conf`

Postgres Pro 11 and below used a configuration file named `recovery.conf` to manage replicas and standbys. Support for this file was removed in Postgres Pro 12. See [the release notes for PostgreSQL 12](#) for details on this change.

On Postgres Pro 12 and above, [archive recovery](#), [streaming replication](#), and [PITR](#) are configured using [normal server configuration parameters](#). These are set in `postgresql.conf` like any other parameter.

The server will not start if a `recovery.conf` exists.

Postgres Pro 15 and below had a setting `promote_trigger_file`, or `trigger_file` before 12. Use `pg_ctl promote` or call `pg_promote()` to promote a standby instead.

The `standby_mode` setting has been removed. A `standby.signal` file in the data directory is used instead. See [Standby Server Operation](#) for details.

## O.2. Default Roles Renamed to Predefined Roles

PostgreSQL 13 and below used the term “Default Roles”. However, as these roles are not able to actually be changed and are installed as part of the system at initialization time, the more appropriate term to use is “Predefined Roles”. See [Section 21.5](#) for current documentation regarding Predefined Roles, and [the release notes for PostgreSQL 14](#) for details on this change.

## O.3. `pg_xlogdump` renamed to `pg_waldump`

Postgres Pro 9.6 and below provided a command named `pg_xlogdump` to read write-ahead-log (WAL) files. This command was renamed to `pg_waldump`, see [pg\\_waldump](#) for documentation of `pg_waldump` and see [the release notes for PostgreSQL 10](#) for details on this change.

## O.4. `pg_resetxlog` renamed to `pg_resetwal`

Postgres Pro 9.6 and below provided a command named `pg_resetxlog` to reset the write-ahead-log (WAL) files. This command was renamed to `pg_resetwal`, see [pg\\_resetwal](#) for documentation of `pg_resetwal` and see [the release notes for PostgreSQL 10](#) for details on this change.

## O.5. `pg_receivexlog` renamed to `pg_receivewal`

Postgres Pro 9.6 and below provided a command named `pg_receivexlog` to fetch write-ahead-log (WAL) files. This command was renamed to `pg_receivewal`, see [pg\\_receivewal](#) for documentation of `pg_receivewal` and see [the release notes for PostgreSQL 10](#) for details on this change.

---

# Bibliography

Selected references and readings for SQL and PostgreSQL.

Some white papers and technical reports from the original POSTGRES development team are available at the University of California, Berkeley, Computer Science Department [web site](#).

## SQL Reference Books

- [bowman01] *The Practical SQL Handbook*. Using SQL Variants. Fourth Edition. Judith Bowman, Sandra Emerson, and Marcy Darnovsky. ISBN 0-201-70309-2. Addison-Wesley Professional. 2001.
- [date97] *A Guide to the SQL Standard*. A user's guide to the standard database language SQL. Fourth Edition. C. J. Date and Hugh Darwen. ISBN 0-201-96426-0. Addison-Wesley. 1997.
- [date04] *An Introduction to Database Systems*. Eighth Edition. C. J. Date. ISBN 0-321-19784-4. Addison-Wesley. 2003.
- [elma04] *Fundamentals of Database Systems*. Fourth Edition. Ramez Elmasri and Shamkant Navathe. ISBN 0-321-12226-7. Addison-Wesley. 2003.
- [melt93] *Understanding the New SQL*. A complete guide. Jim Melton and Alan R. Simon. ISBN 1-55860-245-3. Morgan Kaufmann. 1993.
- [ull88] *Principles of Database and Knowledge-Base Systems*. Classical Database Systems. Jeffrey D. Ullman. Volume 1. Computer Science Press. 1988.
- [sqltr-19075-6] *SQL Technical Report*. Part 6: SQL support for JavaScript Object Notation (JSON). First Edition. 2017.

## PostgreSQL-specific Documentation

- [sim98] *Enhancement of the ANSI SQL Implementation of PostgreSQL*. Stefan Simkovics. Department of Information Systems, Vienna University of Technology. Vienna, Austria. November 29, 1998.
- [yu95] *The Postgres95. User Manual*. A. Yu and J. Chen. University of California. Berkeley, California. Sept. 5, 1995.
- [fong] *The design and implementation of the POSTGRES query optimizer*. Zelaine Fong. University of California, Berkeley, Computer Science Department.

## Proceedings and Articles

- [ports12] “[Serializable Snapshot Isolation in PostgreSQL](#)”. D. Ports and K. Grittner. VLDB Conference, August 2012.
- [berenson95] “[A Critique of ANSI SQL Isolation Levels](#)”. H. Berenson, P. Bernstein, J. Gray, J. Melton, E. O’Neil, and P. O’Neil. ACM-SIGMOD Conference on Management of Data, June 1995.
- [olson93] *Partial indexing in POSTGRES: research project*. Nels Olson. UCB Engin T7.49.1993 O676. University of California. Berkeley, California. 1993.
- [ong90] “A Unified Framework for Version Modeling Using Production Rules in a Database System”. L. Ong and J. Goh. *ERL Technical Memorandum M90/33*. University of California. Berkeley, California. April, 1990.
- [rowe87] “[The POSTGRES data model](#)”. L. Rowe and M. Stonebraker. VLDB Conference, Sept. 1987.



- [seshadri95] "*Generalized Partial Indexes*". P. Seshadri and A. Swami. Eleventh International Conference on Data Engineering, 6–10 March 1995. Cat. No.95CH35724. IEEE Computer Society Press. Los Alamitos, California. 1995. 420–7.
- [ston86] "*The design of POSTGRES*". M. Stonebraker and L. Rowe. ACM-SIGMOD Conference on Management of Data, May 1986.
- [ston87a] "The design of the POSTGRES. rules system". M. Stonebraker, E. Hanson, and C. H. Hong. IEEE Conference on Data Engineering, Feb. 1987.
- [ston87b] "*The design of the POSTGRES storage system*". M. Stonebraker. VLDB Conference, Sept. 1987.
- [ston89] "*A commentary on the POSTGRES rules system*". M. Stonebraker, M. Hearst, and S. Potamianos. *SIGMOD Record* 18(3). Sept. 1989.
- [ston89b] "*The case for partial indexes*". M. Stonebraker. *SIGMOD Record* 18(4). Dec. 1989. 4–11.
- [ston90a] "*The implementation of POSTGRES*". M. Stonebraker, L. A. Rowe, and M. Hirohama. *Transactions on Knowledge and Data Engineering* 2(1). IEEE. March 1990.
- [ston90b] "*On Rules, Procedures, Caching and Views in Database Systems*". M. Stonebraker, A. Jhingran, J. Goh, and S. Potamianos. ACM-SIGMOD Conference on Management of Data, June 1990.
- [ston92] "*An overview of the Sequoia 2000 project*". M. Stonebraker. *Digest of Papers COMPCON Spring 1992*. 1992. 383–388.

---

# Index

## Symbols

\$, 42  
\$libdir, 1174  
\$libdir/plugins, 654, 1835  
\*, 121  
.pgpass, 980  
.pg\_service.conf, 980  
::, 48  
\_PG\_archive\_module\_init, 1463  
\_PG\_init, 1174  
\_PG\_output\_plugin\_init, 1453

## A

abbrev, 272  
ABORT, 1467  
abs, 204  
ACL, 71  
aclcontains, 356  
acldefault, 356  
aclexplode, 357  
aclitem, 75  
aclitemeq, 356  
acos, 208  
acosd, 208  
acosh, 209  
administration tools  
    externally maintained, 2938  
advisory lock, 473  
age, 251, 366  
aggregate function, 21  
    built-in, 334  
    invocation, 44  
    moving aggregate, 1196  
    ordered set, 1199  
    partial aggregation, 1200  
    polymorphic, 1197  
    support functions for, 1201  
    user-defined, 1195  
    variadic, 1197  
alias  
    for table name in query, 21  
    in the FROM clause, 111  
    in the select list, 121  
ALL, 343, 346  
    GROUP BY ALL, 118  
    SELECT ALL, 122  
allow\_alter\_system configuration parameter, 659  
allow\_in\_place\_tablespaces configuration parameter, 665  
allow\_system\_table\_mods configuration parameter, 665  
ALTER AGGREGATE, 1468  
ALTER COLLATION, 1470  
ALTER CONVERSION, 1472  
ALTER DATABASE, 1473  
ALTER DEFAULT PRIVILEGES, 1475  
ALTER DOMAIN, 1478  
ALTER EVENT TRIGGER, 1481  
ALTER EXTENSION, 1482  
ALTER FUNCTION, 1485  
ALTER GROUP, 1488  
ALTER INDEX, 1490  
ALTER LANGUAGE, 1493  
ALTER LARGE OBJECT, 1494  
ALTER MATERIALIZED VIEW, 1495  
ALTER OPERATOR, 1497  
ALTER OPERATOR CLASS, 1499  
ALTER OPERATOR FAMILY, 1500  
ALTER POLICY, 1504  
ALTER PROCEDURE, 1505  
ALTER PROFILE, 1508  
ALTER PUBLICATION, 1510  
ALTER ROLE, 712, 1512  
ALTER ROUTINE, 1516  
ALTER RULE, 1517  
ALTER SCHEMA, 1518  
ALTER SEQUENCE, 1519  
ALTER SERVER, 1523  
ALTER STATISTICS, 1524  
ALTER SUBSCRIPTION, 1525  
ALTER TABLE, 1528  
ALTER TABLESPACE, 1548  
ALTER TEXT SEARCH CONFIGURATION, 1549  
ALTER TEXT SEARCH DICTIONARY, 1551  
ALTER TEXT SEARCH PARSER, 1553  
ALTER TEXT SEARCH TEMPLATE, 1554  
ALTER TRIGGER, 1555  
ALTER TYPE, 1557  
ALTER USER, 1561  
ALTER VIEW, 1562  
ALTER\_REPLICATION\_SLOT, 2319  
amcheck, 2612  
ANALYZE, 752, 1564  
AND (operator), 198  
anonymous code blocks, 1762  
any, 196  
ANY, 337, 343, 346  
anyarray, 196  
anycompatible, 196  
anycompatiblearray, 196  
anycompatiblemultirange, 196  
anycompatiblenonarray, 196  
anycompatiblerange, 196  
anyelement, 196  
anyenum, 196  
anymultirange, 196  
anynonarray, 196  
anyrange, 196  
any\_value, 335  
applicable role, 1105  
application\_name configuration parameter, 633  
aqo, 2617

- arbitrary precision numbers, 135
- Archive Modules, 1463
- archive\_cleanup\_command configuration parameter, 609
- archive\_command configuration parameter, 608
- archive\_library configuration parameter, 608
- archive\_mode configuration parameter, 607
- archive\_timeout configuration parameter, 608
- area, 268
- armor, 2674
- array, 174
  - accessing, 176
  - constant, 175
  - constructor, 50
  - declaration, 174
  - I/O, 181
  - modifying, 178
  - of user-defined type, 1204
  - searching, 180
- ARRAY, 50
  - determination of result type, 407
- array\_agg, 335
- array\_append, 327
- array\_cat, 327
- array\_dims, 327
- array\_fill, 327
- array\_length, 327
- array\_lower, 327
- array\_ndims, 327
- array\_nulls configuration parameter, 657
- array\_position, 327
- array\_positions, 328
- array\_prepend, 328
- array\_remove, 328
- array\_replace, 328
- array\_sample, 328
- array\_shuffle, 328
- array\_to\_json, 296
- array\_to\_string, 328
- array\_to\_tsvector, 275
- array\_upper, 328
- ascii, 212
- asin, 208
- asind, 208
- asinh, 209
- ASSERT
  - in PL/pgSQL, 1316
- assertions
  - in PL/pgSQL, 1316
- asynchronous commit, 847
- AT LOCAL, 260
- AT TIME ZONE, 260
- atan, 208
- atan2, 208
- atan2d, 208
- atand, 208
- atanh, 209
- attach
  - storage parameter, 1633
- authentication\_timeout configuration parameter, 588
- auto-increment (see [serial](#))
- autocommit
  - bulk-loading data, 495
  - psql, 2057
- autoprepared statements, 497
- autoprepares\_for\_protocol configuration parameter, 625
- autoprepares\_limit configuration parameter, 625
- autoprepares\_memory\_limit configuration parameter, 625
- autoprepares\_threshold configuration parameter, 626
- autosummarize storage parameter, 1633
- autovacuum
  - configuration parameters, 643
  - general information, 756
- autovacuum configuration parameter, 643
- autovacuum\_analyze\_scale\_factor
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_analyze\_threshold
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_enabled storage parameter, 1708
- autovacuum\_freeze\_max\_age
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_freeze\_min\_age storage parameter, 1709
- autovacuum\_freeze\_table\_age storage parameter, 1709
- autovacuum\_max\_workers configuration parameter, 643
- autovacuum\_multixact\_freeze\_max\_age
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_multixact\_freeze\_min\_age storage parameter, 1709
- autovacuum\_multixact\_freeze\_table\_age storage parameter, 1710
- autovacuum\_naptime configuration parameter, 643
- autovacuum\_vacuum\_cost\_delay
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_vacuum\_cost\_limit
  - configuration parameter, 645
  - storage parameter, 1709
- autovacuum\_vacuum\_insert\_scale\_factor
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_vacuum\_insert\_threshold
  - configuration parameter, 644
  - storage parameter, 1709
- autovacuum\_vacuum\_scale\_factor
  - configuration parameter, 644

- storage parameter, 1709
- autovacuum\_vacuum\_threshold
  - configuration parameter, 643
  - storage parameter, 1709
- autovacuum\_work\_mem configuration parameter, 594
- auto\_explain, 2639
- auto\_explain.log\_analyze configuration parameter, 2639
- auto\_explain.log\_buffers configuration parameter, 2639
- auto\_explain.log\_format configuration parameter, 2640
- auto\_explain.log\_level configuration parameter, 2640
- auto\_explain.log\_min\_duration configuration parameter, 2639
- auto\_explain.log\_nested\_statements configuration parameter, 2640
- auto\_explain.log\_parameter\_max\_length configuration parameter, 2639
- auto\_explain.log\_settings configuration parameter, 2640
- auto\_explain.log\_timing configuration parameter, 2640
- auto\_explain.log\_triggers configuration parameter, 2640
- auto\_explain.log\_verbose configuration parameter, 2640
- auto\_explain.log\_wal configuration parameter, 2639
- auto\_explain.sample\_rate configuration parameter, 2640
- average, 335
- avg, 335
- B**
- B-Tree (see [index](#))
- backend\_flush\_after configuration parameter, 599
- background freezer, 758
  - configuration parameters, 645
- Background workers, 1444
- backslash escapes, 35
- backslash\_quote configuration parameter, 658
- backtrace\_functions configuration parameter, 665
- backup, 374, 761
- Backup Manifest, 2485
- base type, 1153
- base64 format, 222
- BASE\_BACKUP, 2322
- batch mode
  - in libpq, 952
- BEGIN, 1567
- BETWEEN, 201
- BETWEEN SYMMETRIC, 201
- bgsfreezer\_max\_pages\_per\_rel configuration parameter, 645
- bgsfreezer\_pages\_per\_rel\_to\_clean configuration parameter, 645
- BGWORKER\_BACKEND\_DATABASE\_CONNECTION, 1444
- BGWORKER\_SHMEM\_ACCESS, 1444
- bgwriter\_delay configuration parameter, 598
- bgwriter\_flush\_after configuration parameter, 599
- bgwriter\_lru\_maxpages configuration parameter, 598
- bgwriter\_lru\_multiplier configuration parameter, 599
- bigint, 38, 135
- bigserial, 138
- binary data, 142
  - functions, 219
- binary string
  - concatenation, 219
  - converting to character string, 222
  - length, 221
- bit string
  - constant, 37
  - data type, 160
  - length, 223
- bit strings
  - functions, 223
- bitmap scan, 414, 619
- bit\_and, 335
- bit\_count, 220, 223
- bit\_length, 210, 219, 223
- bit\_or, 335
- bit\_xor, 335
- BLOB (see [large object](#))
- block\_size configuration parameter, 663
- bonjour configuration parameter, 587
- bonjour\_name configuration parameter, 587
- Boolean
  - data type, 153
  - operators (see operators, logical)
- bool\_and, 335
- bool\_or, 335
- booting
  - starting the server during, 561
- bound\_box, 269
- box, 269
- box (data type), 156
- bpchar, 140
- BRIN (see [index](#))
- brin\_desummarize\_range, 385
- brin\_summarize\_new\_values, 384
- brin\_summarize\_range, 384
- broadcast, 272
- BSD Authentication, 708
- btrim, 210, 219
- bt\_index\_check, 2612
- bt\_index\_parent\_check, 2613
- buffering storage parameter, 1632
- bytea, 142
- bytea\_output configuration parameter, 651

**C**

- C, 910, 1008
- C++, 1193
- CALL, 1569
- canceling SQL queries, 957
- cardinality, 328
- CASCADE
  - with DROP, 102
  - foreign key action, 66
- Cascading Replication, 770
- CASE, 323
  - determination of result type, 407
- case sensitivity
  - of SQL commands, 33
- cast
  - I/O conversion, 1599
- cbirt, 204
- ceil, 204
- ceiling, 204
- center, 268
- Certificate, 707
- cfs\_compression\_ratio, 389
- cfs\_compress\_temp\_relations configuration parameter, 661
- cfs\_enable\_gc, 389
- cfs\_estimate, 389
- cfs\_fragmentation, 389
- cfs\_gc configuration parameter, 660
- cfs\_gc\_activity\_processed\_bytes, 389
- cfs\_gc\_activity\_processed\_files, 389
- cfs\_gc\_activity\_processed\_pages, 389
- cfs\_gc\_activity\_scanned\_files, 389
- cfs\_gc\_delay configuration parameter, 661
- cfs\_gc\_lock\_file configuration parameter, 660
- cfs\_gc\_period configuration parameter, 660
- cfs\_gc\_relation, 389
- cfs\_gc\_respond\_time configuration parameter, 661
- cfs\_gc\_threshold configuration parameter, 660
- cfs\_gc\_time\_bloom\_false\_positive configuration parameter, 661
- cfs\_gc\_time\_bloom\_passed configuration parameter, 661
- cfs\_gc\_time\_bloom\_query configuration parameter, 661
- cfs\_gc\_workers configuration parameter, 660
- cfs\_level configuration parameter, 660
- cfs\_log\_verbose configuration parameter, 661
- cfs\_start\_gc, 389
- cfs\_version, 389
- chained transactions, 1579, 1870
  - in PL/pgSQL, 1313
- char, 140
- character, 140
- character set, 653, 664, 740
- character string
  - concatenation, 210
  - constant, 35
  - converting to binary string, 222
  - data types, 140
  - length, 210
  - prefix test, 212
- character varying, 140
- character\_length, 210
- char\_length, 210
- check constraint, 60
- CHECK OPTION, 1749
- checkpoint, 849
- CHECKPOINT, 1570
- checkpoint\_completion\_target configuration parameter, 607
- checkpoint\_flush\_after configuration parameter, 607
- checkpoint\_timeout configuration parameter, 606
- checkpoint\_warning configuration parameter, 607
- checksums, 846
- check\_function\_bodies configuration parameter, 647
- chr, 212
- cid, 193
- cidr, 158
- circle, 157, 269
- client authentication, 689
  - timeout during, 588
- client\_connection\_check\_interval configuration parameter, 588
- client\_encoding configuration parameter, 653
- client\_min\_messages configuration parameter, 645
- clock\_timestamp, 251
- CLOSE, 1571
- cluster
  - of databases (see [database cluster](#))
- CLUSTER, 1572
- cluster initialization, 894
- cluster services, 894
- clusterdb, 1930
- clustering, 770
- cluster\_name configuration parameter, 640
- cmax, 69
- cmin, 68
- COALESCE, 325
- COLLATE, 49
- collation, 731
  - in PL/pgSQL, 1285
  - in SQL functions, 1171
- COLLATION FOR, 361
- colocate\_with storage parameter, 1707
- color, 2974
- column, 16, 56
  - adding, 69
  - removing, 70
  - renaming, 71
  - system column, 68
- column data type
  - changing, 71
- column reference, 42

- col\_description, 364
- comment
  - about database objects, 364
  - in SQL, 40
- COMMENT, 1575
- COMMIT, 1579
- COMMIT PREPARED, 1580
- commit\_delay configuration parameter, 606
- commit\_siblings configuration parameter, 606
- commit\_timestamp\_buffers configuration parameter, 594
- common table expression (see [WITH](#))
- comparison
  - composite type, 346
  - operators, 198
  - row constructor, 346
  - subquery result row, 343
- compiling
  - libpq applications, 985
- composite type, 182, 1153
  - comparison, 346
  - constant, 183
  - constructor, 51
- computed field, 186
- compute\_query\_id configuration parameter, 643
- concat, 213
- concat\_ws, 213
- concurrency, 462
- conditional expression, 323
- configuration
  - of recovery
    - general settings, 608
    - of a standby server, 609
  - of the server, 581
  - of the server
    - functions, 371
- config\_file configuration parameter, 584
- conjunction, 198
- connection service file, 980
- conninfo, 917
- constant, 34
- constraint, 60
  - adding, 70
  - check, 60
  - exclusion, 68
  - foreign key, 65
  - name, 60
  - NOT NULL, 62
  - primary key, 64
  - removing, 70
  - unique, 63
- constraint exclusion, 100, 626
- constraint\_exclusion configuration parameter, 626
- container type, 1153
- CONTINUE
  - in PL/pgSQL, 1301
- continuous archiving, 761
  - in standby, 782
- control file, 1222
- convert, 222
- convert\_from, 222
- convert\_to, 222
- COPY, 18, 1581
  - with libpq, 963
- corr, 338
- correlation, 338
  - in the query planner, 491
- cos, 208
- cosd, 208
- cosh, 209
- cot, 208
- cotd, 208
- count, 335
- covariance
  - population, 338
  - sample, 338
- covar\_pop, 338
- covar\_samp, 338
- covering index, 418
- cpu\_index\_tuple\_cost configuration parameter, 623
- cpu\_operator\_cost configuration parameter, 623
- cpu\_tuple\_cost configuration parameter, 623
- crash\_info configuration parameter, 662
- crash\_info\_dump configuration parameter, 662
- crash\_info\_location configuration parameter, 662
- CREATE ACCESS METHOD, 1591
- CREATE AGGREGATE, 1592
- CREATE CAST, 1599
- CREATE COLLATION, 1603
- CREATE CONVERSION, 1606
- CREATE DATABASE, 723, 1608
- CREATE DOMAIN, 1613
- CREATE EVENT TRIGGER, 1616
- CREATE EXTENSION, 1618
- CREATE FUNCTION, 1620
- CREATE GROUP, 1628
- CREATE INDEX, 1629
- CREATE LANGUAGE, 1637
- CREATE MATERIALIZED VIEW, 1639
- CREATE OPERATOR, 1641
- CREATE OPERATOR CLASS, 1644
- CREATE OPERATOR FAMILY, 1647
- CREATE PACKAGE, 1648
- CREATE POLICY, 1652
- CREATE PROCEDURE, 1657
- CREATE PROFILE, 1661
- CREATE PUBLICATION, 1665
- CREATE ROLE, 710, 1669
- CREATE RULE, 1674
- CREATE SCHEMA, 1677
- CREATE SEQUENCE, 1679
- CREATE SERVER, 1683
- CREATE STATISTICS, 1685
- CREATE SUBSCRIPTION, 1689
- CREATE TABLE, 16, 1694
- CREATE TABLE AS, 1719



- CREATE TABLESPACE, 725, 1722
- CREATE TEXT SEARCH CONFIGURATION, 1724
- CREATE TEXT SEARCH DICTIONARY, 1725
- CREATE TEXT SEARCH PARSER, 1727
- CREATE TEXT SEARCH TEMPLATE, 1729
- CREATE TRANSFORM, 1730
- CREATE TRIGGER, 1732
- CREATE TYPE, 1739
- CREATE USER, 1748
- CREATE VIEW, 1749
- createdb, 11, 1933
- createrole\_self\_grant
  - configuration parameter, 651
- createrole\_self\_grant configuration parameter
  - use in securing functions, 1626
- createuser, 710, 1937
- create\_baseline, 2778
- CREATE\_REPLICATION\_SLOT, 2318
- create\_server, 2771
- cross join, 108
- crypt, 2671
- csn\_commit\_delay configuration parameter, 674
- csn\_lsn\_map\_size configuration parameter, 674
- csn\_max\_shift configuration parameter, 674
- csn\_max\_shift\_error configuration parameter, 674
- csn\_snapshot\_defer\_time configuration parameter, 674
- cstring, 196
- CSV (Comma-Separated Values) format
  - in psql, 2050
- ctid, 69
- CTID, 1258
- CUBE, 118
- cume\_dist, 342
  - hypothetical, 340
- current\_catalog, 352
- current\_database, 352
- current\_date, 251
- current\_logfiles
  - and the log\_destination configuration parameter, 627
  - and the pg\_current\_logfile function, 353
- current\_query, 352
- current\_role, 352
- current\_schema, 352
- current\_schemas, 352
- current\_setting, 371
- current\_time, 251
- current\_timestamp, 251
- current\_user, 352
- currval, 323
- cursor
  - CLOSE, 1571
  - DECLARE, 1755
  - FETCH, 1816
  - in PL/pgSQL, 1308
  - MOVE, 1846
  - showing the query plan, 1810
- cursor\_tuple\_fraction configuration parameter, 626
- custom scan provider
  - handler for, 2378
- D**
- data area (see [database cluster](#))
- data partitioning, 770
- data type, 133
  - base, 1153
  - category, 400
  - composite, 1153
  - constant, 38
  - container, 1153
  - conversion, 399
  - domain, 192
  - enumerated (enum), 154
  - internal organization, 1175
  - numeric, 134
  - polymorphic, 1154
  - type cast, 48
  - user-defined, 1201
- database, 718
  - creating, 11
  - privilege to create, 711
- database activity
  - monitoring, 791
- database cluster, 16, 558
- data\_checksums configuration parameter, 663
- data\_directory configuration parameter, 584
- data\_directory\_mode configuration parameter, 663
- data\_sync\_retry configuration parameter, 661
- date, 144, 145
  - constants, 148
  - current, 262
  - output format, 148
  - (see also [formatting](#))
- DateFormat configuration parameter, 652
- date\_add, 252
- date\_bin, 260
- date\_part, 252, 255
- date\_subtract, 252
- date\_trunc, 252, 259
- dblink, 2642, 2647
  - dblink\_build\_sql\_delete, 2667
  - dblink\_build\_sql\_insert, 2665
  - dblink\_build\_sql\_update, 2668
  - dblink\_cancel\_query, 2663
  - dblink\_close, 2655
  - dblink\_connect, 2643
  - dblink\_connect\_u, 2645
  - dblink\_disconnect, 2646
  - dblink\_error\_message, 2657
  - dblink\_exec, 2650
  - dblink\_fetch, 2653
  - dblink\_get\_connections, 2656
  - dblink\_get\_notify, 2660
  - dblink\_get\_pkey, 2664
  - dblink\_get\_result, 2661

- dblink\_is\_busy, 2659
- dblink\_open, 2652
- dblink\_send\_query, 2658
- deadlock, 472
  - timeout during, 656
- deadlock\_timeout configuration parameter, 656
- DEALLOCATE, 1754
- dearmor, 2674
- debug\_assertions configuration parameter, 663
- debug\_deadlocks configuration parameter, 667
- debug\_discard\_caches configuration parameter, 665
- debug\_io\_direct configuration parameter, 665
- debug\_logical\_replication\_streaming configuration parameter, 669
- debug\_parallel\_query configuration parameter, 666
- debug\_pretty\_print configuration parameter, 633
- debug\_print\_parse configuration parameter, 633
- debug\_print\_plan configuration parameter, 633
- debug\_print\_rewritten configuration parameter, 633
- decimal (see [numeric](#))
- DECLARE, 1755
- decode, 222
- decode\_bytea
  - in PL/Perl, 1358
- decrypt, 2676
- decrypt\_iv, 2676
- deduplicate\_items storage parameter, 1632
- default value, 57
  - changing, 70
- default-roles, 2975
- default\_statistics\_target configuration parameter, 626
- default\_tablespace configuration parameter, 646
- default\_table\_access\_method configuration parameter, 646
- default\_text\_search\_config configuration parameter, 654
- default\_toast\_compression configuration parameter, 647
- default\_transaction\_deferrable configuration parameter, 647
- default\_transaction\_isolation configuration parameter, 647
- default\_transaction\_read\_only configuration parameter, 647
- deferrable transaction, 648
  - setting, 1908
  - setting default, 647
- degrees, 204
- delay, 263
- DELETE, 23, 105, 1758
  - RETURNING, 106
- deleting, 105
- dense\_rank, 342
  - hypothetical, 340
- diagonal, 268
- diameter, 268
- digest, 2670
- dirty read, 462
- disable\_server, 2772
- DISCARD, 1761
- disjunction, 198
- disk drive, 853
- disk space, 751
- disk usage, 842
- DISTINCT, 19
  - GROUP BY DISTINCT, 118
  - SELECT DISTINCT, 122
- distributed\_by storage parameter, 1707
- div, 204
- DO, 1762
- document
  - text search, 427
- dollar quoting, 36
- domain, 192
- double precision, 137
- DROP ACCESS METHOD, 1763
- DROP AGGREGATE, 1764
- DROP CAST, 1766
- DROP COLLATION, 1767
- DROP CONVERSION, 1768
- DROP DATABASE, 725, 1769
- DROP DOMAIN, 1770
- DROP EVENT TRIGGER, 1771
- DROP EXTENSION, 1772
- DROP FUNCTION, 1773
- DROP GROUP, 1775
- DROP INDEX, 1776
- DROP LANGUAGE, 1777
- DROP MATERIALIZED VIEW, 1778
- DROP OPERATOR, 1779
- DROP OPERATOR CLASS, 1780
- DROP OPERATOR FAMILY, 1781
- DROP OWNED, 1782
- DROP PACKAGE, 1783
- DROP POLICY, 1784
- DROP PROCEDURE, 1785
- DROP PROFILE, 1787
- DROP PUBLICATION, 1788
- DROP ROLE, 710, 1789
- DROP ROUTINE, 1790
- DROP RULE, 1791
- DROP SCHEMA, 1792
- DROP SEQUENCE, 1793
- DROP SERVER, 1794
- DROP STATISTICS, 1795
- DROP SUBSCRIPTION, 1796
- DROP TABLE, 18, 1797
- DROP TABLESPACE, 1798
- DROP TEXT SEARCH CONFIGURATION, 1799
- DROP TEXT SEARCH DICTIONARY, 1800
- DROP TEXT SEARCH PARSER, 1801
- DROP TEXT SEARCH TEMPLATE, 1802
- DROP TRANSFORM, 1803



DROP TRIGGER, 1804  
DROP TYPE, 1805  
DROP USER, 1806  
DROP VIEW, 1807  
dropdb, 1942  
dropuser, 710, 1945  
drop\_baseline, 2778  
DROP\_REPLICATION\_SLOT, 2322  
drop\_server, 2772  
DTD, 164

duplicate, 19  
duplicates, 122  
dynamic loading, 655, 1174  
dynamic\_library\_path, 1174  
dynamic\_library\_path configuration parameter, 655  
dynamic\_shared\_memory\_type configuration parameter, 596

## E

ECPG, 1008  
ecpg, 1947  
EDITION, 2099  
effective\_cache\_size configuration parameter, 624  
effective\_io\_concurrency configuration parameter, 599  
elog  
    in PL/Perl, 1358  
    in PL/Python, 1377  
    in PL/Tcl, 1344  
embedded SQL  
    in C, 1008  
enabled role, 1122  
enable\_alternative\_sorting\_cost\_model configuration parameter, 620  
enable\_any\_to\_lateral\_transformation configuration parameter, 619  
enable\_appendopath configuration parameter, 627  
enable\_async\_append configuration parameter, 619  
enable\_async\_merge\_append configuration parameter, 673  
enable\_background\_freezer configuration parameter, 645  
enable\_bitmapscan configuration parameter, 619  
enable\_compound\_index\_stats configuration parameter, 621  
ENABLE\_CRASH\_INFO, 2099  
enable\_csn\_snapshot configuration parameter, 657  
enable\_extra\_transformations configuration parameter, 620  
enable\_gathermerge configuration parameter, 620  
enable\_group\_by\_reordering configuration parameter, 620  
enable\_hashagg configuration parameter, 620  
enable\_hashjoin configuration parameter, 620  
enable\_incremental\_sort configuration parameter, 620  
enable\_indexonlyscan configuration parameter, 620

enable\_indexscan configuration parameter, 620  
enable\_large\_mem\_buffers configuration parameter, 658  
enable\_material configuration parameter, 620  
enable\_memoize configuration parameter, 620  
enable\_mergejoin configuration parameter, 620  
enable\_merge\_append configuration parameter, 673  
enable\_nestloop configuration parameter, 620  
ENABLE\_NLS, 2099  
enable\_non\_equivalence\_filters configuration parameter, 682  
enable\_parallel\_append configuration parameter, 621  
enable\_parallel\_hash configuration parameter, 621  
enable\_parallel\_temptables configuration parameter, 621  
enable\_partitionwise\_aggregate configuration parameter, 621  
enable\_partitionwise\_join configuration parameter, 621  
enable\_partition\_pruning configuration parameter, 621  
enable\_partition\_pruning\_extra configuration parameter, 675  
ENABLE\_PGPRO\_TUNE, 2099  
enable\_presorted\_aggregate configuration parameter, 621  
enable\_self\_join\_removal configuration parameter, 622  
enable\_seqscan configuration parameter, 622  
enable\_server, 2772  
enable\_sort configuration parameter, 622  
enable\_sql\_func\_custom\_plans configuration parameter, 673  
enable\_temp\_memory\_catalog configuration parameter, 596  
enable\_tidscan configuration parameter, 622  
encode, 222  
encode\_array\_constructor  
    in PL/Perl, 1358  
encode\_array\_literal  
    in PL/Perl, 1358  
encode\_bytea  
    in PL/Perl, 1358  
encode\_typed\_literal  
    in PL/Perl, 1358  
encrypt, 2676  
encryption, 574  
    for specific columns, 2670  
encrypt\_iv, 2676  
END, 1808  
enumerated types, 154  
enum\_first, 264  
enum\_last, 264  
enum\_range, 264  
environment variable, 978  
ephemeral named relation

- registering with SPI, 1416, 1418
- unregistering from SPI, 1417
- erf, 204
- erfc, 205
- error codes
  - libpq, 939
  - list of, 2505
- error message
  - in PGcancelConn, 959
  - in PGconn, 931
- escape format, 222
- escape string syntax, 35
- escape\_string\_warning configuration parameter, 658
- escaping strings
  - in libpq, 945
- event log
  - event log, 580
- event trigger, 1243
  - in C, 1247
  - in PL/Tcl, 1346
- event\_source configuration parameter, 631
- event\_trigger, 196
- event\_triggers
  - configuration parameter, 652
- every, 335
- EXCEPT, 122
- exceptions
  - in PL/pgSQL, 1304
  - in PL/Tcl, 1347
- exclusion constraint, 68
- EXECUTE, 1809
- EXISTS, 343
- EXIT
  - in PL/pgSQL, 1300
- exit\_on\_error configuration parameter, 661
- exp, 205
- EXPLAIN, 477, 1810
- export\_data, 2779
- expression
  - order of evaluation, 52
  - syntax, 41
- extending SQL, 1153
- extension, 1221
  - externally maintained, 2938
- external\_pid\_file configuration parameter, 585
- extract, 252, 255
- extra\_float\_digits configuration parameter, 652

## F

- factorial, 205
- failover, 770
- false, 153
- family, 272
- fast path, 961
- fastupdate storage parameter, 1632
- fault tolerance, 894
- fdw\_handler, 196

- FETCH, 1816
- field
  - computed, 186
- field selection, 43
- file system mount points, 560
- fillfactor storage parameter, 1632, 1708
- FILTER, 44
- first\_value, 342
- float4 (see [real](#))
- float8 (see [double precision](#))
- floating point, 137
- floating-point
  - display, 652
- floor, 205
- foreign data, 101
- foreign data wrapper
  - handler for, 2356
- foreign key, 25, 65
  - self-referential, 66
- foreign table, 101
- foreign\_analyze\_interval configuration parameter, 674
- foreign\_join\_fast\_path configuration parameter, 674
- format, 213, 217
  - use in PL/pgSQL, 1292
- formatting, 242
- format\_type, 358
- Free Space Map, 2466
- FreeBSD
  - IPC configuration, 565
  - shared library, 1181
  - start script, 561
- from\_collapse\_limit configuration parameter, 626
- FSM (see [Free Space Map](#))
- fsync configuration parameter, 602
- full text search, 426
  - data types, 161
  - functions and operators, 161
- full\_page\_writes configuration parameter, 604
- function, 198
  - default values for arguments, 1164
  - in the FROM clause, 112
  - internal, 1174
  - invocation, 43
  - mixed notation, 55
  - named argument, 1158
  - named notation, 54
  - output parameter, 1162
  - polymorphic, 1154
  - positional notation, 54
  - RETURNS TABLE, 1169
  - statistics, 395
  - type resolution in an invocation, 403
  - user-defined, 1156
    - in C, 1174
    - in SQL, 1157
  - variadic, 1164

with SETOF, 1166  
functional dependency, 117

## G

gcd, 205  
generated column, 59, 1703  
    in triggers, 1235  
generate\_series, 348  
generate\_subscripts, 350  
generic\_plan\_fuzz\_factor configuration parameter, 624  
genetic query optimization, 624  
gen\_random\_bytes, 2677  
gen\_random\_uuid, 279, 2677  
gen\_salt, 2671  
GEQO (see [genetic query optimization](#))  
geqo configuration parameter, 624  
geqo\_effort configuration parameter, 625  
geqo\_generations configuration parameter, 625  
geqo\_pool\_size configuration parameter, 625  
geqo\_seed configuration parameter, 625  
geqo\_selection\_bias configuration parameter, 625  
geqo\_threshold configuration parameter, 625  
get\_bit, 220, 224  
get\_byte, 220  
get\_current\_ts\_config, 275  
get\_diffreport, 2781  
get\_report, 2780  
get\_report\_latest, 2780  
GIN (see [index](#))  
gin\_clean\_pending\_list, 385  
gin\_fuzzy\_search\_limit configuration parameter, 656  
gin\_pending\_list\_limit  
    configuration parameter, 651  
    storage parameter, 1633  
GiST (see [index](#))  
global data  
    in PL/Python, 1371  
    in PL/Tcl, 1342  
global storage parameter, 1542, 1707  
Global Views for Statistics, 2292  
Global Views for System Catalog and System Views, 2297  
GRANT, 71, 1820  
GREATEST, 326  
    determination of result type, 407  
Gregorian calendar, 2519  
GROUP BY, 22, 116  
grouping, 116  
GROUPING, 341  
GROUPING SETS, 118  
gssapi, 579  
GSSAPI, 700  
    with libpq, 923  
gss\_accept\_delegation configuration parameter, 589  
GUID, 163

## H

hash (see [index](#))  
hash\_mem\_multiplier configuration parameter, 594  
has\_any\_column\_privilege, 355  
has\_column\_privilege, 355  
has\_database\_privilege, 355  
has\_foreign\_data\_wrapper\_privilege, 355  
has\_function\_privilege, 355  
has\_language\_privilege, 355  
has\_parameter\_privilege, 355  
has\_schema\_privilege, 355  
has\_sequence\_privilege, 355  
has\_server\_privilege, 355  
has\_tablespace\_privilege, 355  
has\_table\_privilege, 355  
has\_type\_privilege, 356  
HAVING, 22, 117  
hba\_file configuration parameter, 585  
height, 268  
hex format, 222  
hierarchical database, 16  
high availability, 770  
history  
    of PostgreSQL, xxv  
hmac, 2670  
host, 273  
host name, 919  
hostmask, 273  
hot standby, 770  
hot\_standby configuration parameter, 616  
hot\_standby\_feedback configuration parameter, 617  
huge\_pages configuration parameter, 592  
huge\_pages\_status configuration parameter, 663  
huge\_page\_size configuration parameter, 592  
hypothetical-set aggregate  
    built-in, 340

## I

ICU, 729, 733, 1603, 1610  
icu\_unicode\_version, 370  
icu\_validation\_level configuration parameter, 653  
ident, 702  
identifier  
    length, 33  
    syntax of, 33  
IDENTIFY\_SYSTEM, 2317  
identity column, 58  
ident\_file configuration parameter, 585  
idle\_in\_transaction\_session\_timeout configuration parameter, 649  
idle\_session\_timeout configuration parameter, 649  
IFNULL, 325  
ignore\_checksum\_failure configuration parameter, 668  
ignore\_invalid\_pages configuration parameter, 668

- ignore\_system\_indexes configuration parameter, 666
  - IMMUTABLE, 1172
  - import\_data, 2779
  - IN, 343, 346
  - INCLUDE
    - in index definitions, 419
  - include
    - in configuration file, 583
  - include\_dir
    - in configuration file, 583
  - include\_if\_exists
    - in configuration file, 583
  - index, 410
    - and ORDER BY, 414
    - B-Tree, 411, 2406
    - BRIN, 412, 2444
    - building concurrently, 1633
    - combining multiple indexes, 414
    - covering, 418
    - examining usage, 422
    - on expressions, 415
    - for user-defined data type, 1209
    - GIN, 412, 2438
      - text search, 457
    - GiST, 412, 2412
      - text search, 457
    - hash, 411
    - Hash, 2456
    - index-only scans, 418
    - locks, 476
    - multicolumn, 413
    - partial, 416
    - rebuilding concurrently, 1860
    - RUM
      - text search, 457
    - SP-GiST, 412, 2427
    - unique, 415
  - Index Access Method, 2388
  - index scan, 620
  - index-only scan, 418
  - indexam
    - Index Access Method, 2388
  - index\_am\_handler, 196
  - inet (data type), 158
  - inet\_client\_addr, 352
  - inet\_client\_port, 352
  - inet\_merge, 273
  - inet\_same\_family, 273
  - inet\_server\_addr, 352
  - inet\_server\_port, 352
  - infinity
    - floating point, 138
    - numeric (data type), 136
  - information schema, 1104
  - inheritance, 30, 86
  - initcap, 213
  - initdb, 558, 2079
  - Initialization Fork, 2467
  - initplan, 484
  - input function, 1201
  - INSERT, 18, 104, 1826
    - RETURNING, 106
  - inserting, 104
  - instr function, 1338
  - int2 (see [smallint](#))
  - int4 (see [integer](#))
  - int8 (see [bigint](#))
  - integer, 38, 135
  - integer\_datetimes configuration parameter, 663
  - interfaces
    - externally maintained, 2938
  - internal, 196
  - INTERSECT, 122
  - interval, 144, 150
    - output format, 152
    - (see also [formatting](#))
  - IntervalStyle configuration parameter, 652
  - inverse distribution, 339
  - in\_hot\_standby configuration parameter, 663
  - in\_range support functions, 2407
  - io\_combine\_limit configuration parameter, 600
  - IS DISTINCT FROM, 201, 346
  - IS DOCUMENT, 284
  - IS FALSE, 202
  - IS JSON, 298
  - IS NOT DISTINCT FROM, 201, 346
  - IS NOT DOCUMENT, 284
  - IS NOT FALSE, 202
  - IS NOT NULL, 201
  - IS NOT TRUE, 202
  - IS NOT UNKNOWN, 202
  - IS NULL, 201, 659
  - IS TRUE, 202
  - IS UNKNOWN, 202
  - isclosed, 268
  - isempty, 333, 334
  - isfinite, 252
  - ISNULL, 201
  - isopen, 268
  - is\_array\_ref
    - in PL/Perl, 1359
- ## J
- JIT, 876
  - jit configuration parameter, 627
  - jit\_above\_cost configuration parameter, 624
  - jit\_debugging\_support configuration parameter, 668
  - jit\_dump\_bitcode configuration parameter, 669
  - jit\_expressions configuration parameter, 669
  - jit\_inline\_above\_cost configuration parameter, 624
  - jit\_optimize\_above\_cost configuration parameter, 624
  - jit\_profiling\_support configuration parameter, 669
  - jit\_provider configuration parameter, 655

- jit\_tuple\_deforming configuration parameter, 669
  - join, 20, 108
    - controlling the order, 494
    - cross, 108
    - left, 109
    - natural, 109
    - outer, 21, 108
    - right, 109
    - self, 21
  - join\_collapse\_limit configuration parameter, 627
  - JSON, 165
    - functions and operators, 292
  - json constructor, 297
  - JSONB, 165
  - jsonb
    - containment, 167
    - existence, 168
    - indexes on, 169
  - jsonb\_agg, 336
  - jsonb\_agg\_strict, 336
  - jsonb\_array\_elements, 298
  - jsonb\_array\_elements\_text, 299
  - jsonb\_array\_length, 299
  - jsonb\_build\_array, 296
  - jsonb\_build\_object, 296
  - jsonb\_each, 299
  - jsonb\_each\_text, 299
  - jsonb\_extract\_path, 299
  - jsonb\_extract\_path\_text, 299
  - jsonb\_insert, 302
  - jsonb\_object, 297
  - jsonb\_object\_agg, 336
  - jsonb\_object\_agg\_strict, 336
  - jsonb\_object\_agg\_unique, 336
  - jsonb\_object\_agg\_unique\_strict, 336
  - jsonb\_object\_keys, 299
  - jsonb\_path\_exists, 302
  - jsonb\_path\_exists\_tz, 303
  - jsonb\_path\_match, 303
  - jsonb\_path\_match\_tz, 303
  - jsonb\_path\_query, 303
  - jsonb\_path\_query\_array, 303
  - jsonb\_path\_query\_array\_tz, 303
  - jsonb\_path\_query\_first, 303
  - jsonb\_path\_query\_first\_tz, 304
  - jsonb\_path\_query\_tz, 303
  - jsonb\_populate\_record, 300
  - jsonb\_populate\_recordset, 301
  - jsonb\_populate\_record\_valid, 300
  - jsonb\_pretty, 304
  - jsonb\_set, 302
  - jsonb\_set\_lax, 302
  - jsonb\_strip\_nulls, 302
  - jsonb\_to\_record, 301
  - jsonb\_to\_recordset, 301
  - jsonb\_to\_tsvector, 277
  - jsonb\_typeof, 304
  - jsonpath, 172
  - json\_agg, 336
  - json\_agg\_strict, 336
  - json\_array, 296
  - json\_arrayagg, 336
  - json\_array\_elements, 298
  - json\_array\_elements\_text, 299
  - json\_array\_length, 299
  - json\_build\_array, 296
  - json\_build\_object, 296
  - json\_each, 299
  - json\_each\_text, 299
  - json\_exists, 314
  - json\_extract\_path, 299
  - json\_extract\_path\_text, 299
  - json\_object, 297, 297
  - json\_objectagg, 336
  - json\_object\_agg, 336
  - json\_object\_agg\_strict, 336
  - json\_object\_agg\_unique, 336
  - json\_object\_agg\_unique\_strict, 336
  - json\_object\_keys, 299
  - json\_populate\_record, 300
  - json\_populate\_recordset, 301
  - json\_query, 315
  - json\_scalar, 297
  - json\_strip\_nulls, 302
  - json\_table, 316
  - json\_to\_record, 301
  - json\_to\_recordset, 301
  - json\_to\_tsvector, 277
  - json\_typeof, 304
  - json\_value, 316
  - Julian date, 2520
  - Just-In-Time compilation (see [JIT](#))
  - justify\_days, 253
  - justify\_hours, 253
  - justify\_interval, 253
- K**
- keep\_baseline, 2778
  - key word
    - list of, 2522
    - syntax of, 33
  - krb\_caseins\_users configuration parameter, 589
  - krb\_server\_keyfile configuration parameter, 588
- L**
- label (see [alias](#))
  - lag, 342
  - language\_handler, 196
  - large object, 997
  - lastval, 323
  - last\_value, 342
  - LATERAL
    - in the FROM clause, 114
  - lcm, 205
  - lc\_collate configuration parameter, 663
  - lcctype configuration parameter, 663

- lc\_messages configuration parameter, 653
- lc\_monetary configuration parameter, 653
- lc\_numeric configuration parameter, 653
- lc\_time configuration parameter, 653
- LDAP, 703
- LDAP connection parameter lookup, 981
- lead, 342
- LEAST, 326
  - determination of result type, 407
- left, 213
- left join, 109
- length, 213, 221, 223, 268, 275
  - of a binary string (see [binary strings](#), [length](#))
  - of a character string (see [character string](#), [length](#))
- length(tsvector), 438
- libedit
  - in psql, 2063
- libpq, 910
  - chunked mode, 956
  - pipeline mode, 952
  - single-row mode, 956
- libpq-fe.h, 910, 928
- libpq-int.h, 928
- library initialization function, 1174
- LIKE, 225
  - and locales, 728
- LIKE\_REGEX, 240
  - in SQL/JSON, 314
- LIMIT, 124
- line, 156, 270
- line segment, 156
- linear regression, 338
- Linux
  - IPC configuration, 566
  - shared library, 1181
  - start script, 561
- LISTEN, 1833
- listen\_addresses configuration parameter, 585
- ln, 205
- LOAD, 1835
- load balancing, 770
- locale, 559, 727
- localtime, 253
- localtimestamp, 253
- local\_preload\_libraries configuration parameter, 654
- lock, 468
  - advisory, 473
  - monitoring, 834
- LOCK, 468, 1836
- lock\_timeout configuration parameter, 649
- log, 205
- log shipping, 770
- log10, 205
- log2\_num\_lock\_partitions configuration parameter, 657
- Logging
  - current\_logfiles file and the pg\_current\_logfile function, 353
  - pg\_current\_logfile function, 353
- logging\_collector configuration parameter, 628
- Logical Decoding, 1447, 1450
- logical\_decoding\_work\_mem configuration parameter, 594
- login privilege, 711
- log\_autovacuum\_min\_duration
  - configuration parameter, 633
  - storage parameter, 1710
- log\_btree\_build\_stats configuration parameter, 667
- log\_checkpoints configuration parameter, 633
- log\_connections configuration parameter, 634
- log\_destination configuration parameter, 627
- log\_directory configuration parameter, 629
- log\_disconnections configuration parameter, 634
- log\_duration configuration parameter, 634
- log\_error\_verbosity configuration parameter, 634
- log\_executor\_stats configuration parameter, 643
- log\_filename configuration parameter, 629
- log\_file\_mode configuration parameter, 629
- log\_hostname configuration parameter, 634
- log\_line\_prefix configuration parameter, 634
- log\_lock\_waits configuration parameter, 636
- log\_min\_duration\_sample configuration parameter, 631
- log\_min\_duration\_statement configuration parameter, 631
- log\_min\_error\_statement configuration parameter, 631
- log\_min\_messages configuration parameter, 631
- log\_next\_xid\_assign\_threshold configuration parameter, 637
- log\_parameter\_max\_length configuration parameter, 636
- log\_parameter\_max\_length\_on\_error configuration parameter, 637
- log\_parser\_stats configuration parameter, 643
- log\_planner\_stats configuration parameter, 643
- log\_recovery\_conflict\_waits configuration parameter, 636
- log\_replication\_commands configuration parameter, 637
- log\_rotation\_age configuration parameter, 629
- log\_rotation\_size configuration parameter, 629
- log\_startup\_progress\_interval configuration parameter, 632
- log\_statement configuration parameter, 637
- log\_statement\_sample\_rate configuration parameter, 632
- log\_statement\_stats configuration parameter, 643
- log\_temp\_files configuration parameter, 638
- log\_timezone configuration parameter, 638
- log\_transaction\_sample\_rate configuration parameter, 632
- log\_truncate\_on\_rotation configuration parameter, 630



- looks\_like\_number
  - in PL/Perl, 1358
- loop
  - in PL/pgSQL, 1300
- lower, 210, 333, 333
  - and locales, 728
- lower\_inc, 333, 334
- lower\_inf, 333, 334
- lo\_close, 1000
- lo\_compat\_privileges configuration parameter, 658
- lo\_creat, 998, 1001
- lo\_create, 998
- lo\_export, 998, 1001
- lo\_from\_bytea, 1001
- lo\_get, 1001
- lo\_import, 998, 1001
- lo\_import\_with\_oid, 998
- lo\_lseek, 999
- lo\_lseek64, 999
- lo\_open, 998
- lo\_put, 1001
- lo\_read, 999
- lo\_tell, 1000
- lo\_tell64, 1000
- lo\_truncate, 1000
- lo\_truncate64, 1000
- lo\_unlink, 1000, 1001
- lo\_write, 999
- lpad, 211
- lseg, 156, 270
- LSN, 852
- ltrim, 211, 219
- lwlock\_shared\_limit configuration parameter, 657

## M

- MAC address (see `macaddr`)
- MAC address (EUI-64 format) (see `macaddr`)
- `macaddr` (data type), 159
- `macaddr8` (data type), 159
- `macaddr8_set7bit`, 274
- macOS
  - IPC configuration, 566
  - shared library, 1182
- magic block, 1174
- maintenance, 750
- `maintenance_io_concurrency` configuration parameter, 600
- `maintenance_work_mem` configuration parameter, 594
- `makeaclitem`, 357
- `make_date`, 253
- `make_interval`, 253
- `make_time`, 253
- `make_timestamp`, 253
- `make_timestamptz`, 254
- `masklen`, 273
- materialized view
  - implementation through rules, 1260
- materialized views, 2269
- max, 336
- `max_autonomous_transactions` configuration parameter, 593
- `max_backend_memory` configuration parameter, 597
- `max_connections` configuration parameter, 585
- `max_files_per_process` configuration parameter, 597
- `max_function_args` configuration parameter, 663
- `max_identifier_length` configuration parameter, 663
- `max_index_keys` configuration parameter, 663
- `max_locks_per_transaction` configuration parameter, 656
- `max_logical_replication_workers` configuration parameter, 619
- `max_notify_queue_pages` configuration parameter, 596
- `max_parallel_apply_workers_per_subscription` configuration parameter, 619
- `max_parallel_autovacuum_workers` configuration parameter, 645
- `max_parallel_maintenance_workers` configuration parameter, 600
- `max_parallel_workers` configuration parameter, 601
- `max_parallel_workers_per_gather` configuration parameter, 600
- `max_pred_locks_per_page` configuration parameter, 657
- `max_pred_locks_per_relation` configuration parameter, 657
- `max_pred_locks_per_transaction` configuration parameter, 656
- `max_prepared_transactions` configuration parameter, 593
- `max_replication_slots` configuration parameter
  - in a sending server, 613
  - in a subscriber, 618
- `max_slot_wal_keep_size` configuration parameter, 613
- `max_stack_depth` configuration parameter, 595
- `max_standby_archive_delay` configuration parameter, 616
- `max_standby_streaming_delay` configuration parameter, 616
- `max_sync_workers_per_subscription` configuration parameter, 619
- `max_wal_senders` configuration parameter, 612
- `max_wal_size` configuration parameter, 607
- `max_worker_processes` configuration parameter, 600
- md5, 213, 221
- MD5, 699
- median, 45
  - (see also [percentile](#))
- MEMMB, 2099

- memory context
    - in SPI, 1427
  - memory overcommit, 569
  - MERGE, 1839
    - RETURNING, 106, 343
  - merge\_action, 343
  - min, 337
  - min\_dynamic\_shared\_memory configuration parameter, 596
  - min\_parallel\_index\_scan\_size configuration parameter, 624
  - min\_parallel\_table\_scan\_size configuration parameter, 623
  - min\_scale, 205
  - min\_wal\_size configuration parameter, 607
  - mod, 206
  - mode
    - statistical, 339
  - monitoring
    - database activity, 791
  - MOVE, 1846
  - moving-aggregate mode, 1196
  - multirange (function), 334
  - multirange type, 188
  - Multiversion Concurrency Control, 462
  - MultiXactId, 755
  - multixact\_member\_buffers configuration parameter, 595
  - multixact\_offset\_buffers configuration parameter, 595
  - MVCC, 462
  - MVER, 2099
  - mxid\_age, 366
- N**
- name
    - qualified, 82
    - syntax of, 33
    - unqualified, 83
  - NaN (see [not a number](#))
  - natural join, 109
  - NCPU, 2099
  - negation, 198
  - NetBSD
    - IPC configuration, 566
    - shared library, 1182
    - start script, 562
  - netmask, 273
  - network, 273
    - data types, 157
  - NETWORK EXPLAIN ANALYZE parameter, 504
  - nextval, 322
  - NFS, 560
  - nonblocking connection, 912, 914, 948
  - nonrepeatable read, 462
  - normalize, 211
  - normalized, 210
  - NOT (operator), 198
  - not a number
    - floating point, 138
    - numeric (data type), 136
  - NOT IN, 343, 346
  - not-null constraint, 62
  - notation
    - functions, 53
  - notice processing
    - in libpq, 971
  - notice processor, 972
  - notice receiver, 972
  - NOTIFY, 1848
    - in libpq, 962
  - notify\_buffers configuration parameter, 595
  - NOTNULL, 201
  - now, 254
  - npoints, 269
  - nth\_value, 342
  - ntile, 342
  - null value
    - with check constraints, 62
    - comparing, 201
    - default value, 57
    - in DISTINCT, 122
    - in libpq, 944
    - in PL/Perl, 1351
    - in PL/Python, 1367
    - with unique constraints, 64
  - NULLIF, 325
  - nul\_byte\_replacement\_on\_import configuration parameter, 658
  - number
    - constant, 37
  - numeric, 38
  - numeric (data type), 135
  - numnode, 275, 439
  - num\_nonnulls, 202
  - num\_nulls, 202
  - num\_parts storage parameter, 1707
  - NVL, 325
- O**
- object identifier
    - data type, 193
  - object-oriented database, 16
  - obj\_description, 365
  - OCCURRENCES\_REGEX, 240
  - octet\_length, 211, 211, 220, 224
  - OFFSET, 124
  - oid, 193
  - OID
    - in libpq, 945
  - oid2name, 2832
  - ON CONFLICT, 1826
  - ONLY, 108
  - OOM, 569
  - OpenBSD
    - IPC configuration, 566



- shared library, 1182
- start script, 561
- operator, 198
  - invocation, 43
  - logical, 198
  - precedence, 40
  - syntax, 39
  - type resolution in an invocation, 400
  - user-defined, 1205
- operator class, 421, 1209
- operator family, 421, 1216
- optimization information
  - for functions, 1194
  - for operators, 1206
- optimize\_correlated\_subqueries configuration parameter, 674
- optimize\_row\_in\_expr configuration parameter, 682
- OR (operator), 198
- Oracle
  - porting from PL/SQL to PL/pgSQL, 1331
- ORDER BY, 19, 123
  - and locales, 728
- ordered-set aggregate, 44
  - built-in, 339
- ordering operator, 1219
- order\_by\_attach
  - storage parameter, 1633
- ordinality, 351
- outer join, 108
- output function, 1201
- OVER clause, 46
- overcommit, 569
- OVERLAPS, 254
- overlay, 211, 220, 224
- overloading
  - functions, 1171
  - operators, 1205
- owner, 71

## P

- packages
  - in PL/pgSQL, 1324
- pages\_per\_range storage parameter, 1633
- page\_repair configuration parameter, 613
- pallo, 1181
- PAM, 708
- parallel query, 510
- parallel\_autovacuum\_workers storage parameter, 1708
- parallel\_leader\_participation configuration parameter, 601
- parallel\_setup\_cost configuration parameter, 623
- parallel\_tuple\_cost configuration parameter, 623
- parallel\_workers storage parameter, 1708
- parameter
  - syntax, 42
- parenthesis, 42
- parse\_ident, 213
- partition pruning, 99
- partitioned table, 89
- partitioning, 89
- partition\_bounds storage parameter, 1707
- partition\_by storage parameter, 1707
- password, 711
  - authentication, 699
  - of the superuser, 559
- password file, 980
- password\_encryption configuration parameter, 588
- path, 270
  - for schemas, 645
- path (data type), 157
- pattern matching, 224
- patterns
  - in psql and pg\_dump, 2055
- pclose, 269
- peer, 703
- percentile
  - continuous, 340
  - discrete, 340
- percent\_rank, 342
  - hypothetical, 340
- performance, 477
- Perl, 1350
- permission (see [privilege](#))
- pfree, 1181
- pg-setup, 2110
- pg-wrapper, 2023
- PGAPPNAME, 979
- pgbench, 1954
- pgbouncer, 2905
- PGcancel, 960
- PGcancelConn, 957
- PGCHANNELBINDING, 978
- PGCLIENTENCODING, 979
- PGconn, 910
- PGCONNECT\_TIMEOUT, 979
- pgcrypto, 2670
- PGDATA, 559
- PGDATABASE, 978
- PGDATESTYLE, 979
- PGEventProc, 974
- PGGEQO, 979
- PGGSSDELEGATION, 979
- PGGSSENCMODE, 979
- PGGSSLIB, 979
- PGHOST, 978
- PGHOSTADDR, 978
- PGKRBSRVNAME, 979
- PGLOADBALANCEHOSTS, 979
- PGLOCALEDIR, 979
- PGOPTIONS, 979
- PGPASSFILE, 978
- PGPASSWORD, 978
- PGPORT, 978
- pgpro\_build, 370
- pgpro\_build configuration parameter, 664

- pgpro\_controldata, 2901
- pgpro\_edition, 370
- pgpro\_edition configuration parameter, 664
- pgpro\_pwr, 2766
- pgpro\_pwr.max configuration parameter, 2770
- pgpro\_pwr.max\_query\_length configuration parameter, 2771
- pgpro\_pwr.max\_sample\_age configuration parameter, 2770
- pgpro\_pwr.relsizes\_collect\_mode configuration parameter, 2771
- pgpro\_pwr.statements\_reset configuration parameter, 2771
- pgpro\_pwr.track\_sample\_timings configuration parameter, 2771
- pgpro\_stats, 2692
- pgpro\_stats.enable\_inval\_msgs\_counters configuration parameter, 2733
- pgpro\_stats\_create\_pg\_stat\_kcache\_compatible\_views
  - function, 2719
- pgpro\_stats\_create\_pg\_stat\_statements\_compatible\_views
  - function, 2719
- pgpro\_stats\_get\_archiver
  - function, 2719
- pgpro\_stats\_info, 2718
- pgpro\_stats\_inval\_status
  - function, 2732
- pgpro\_stats\_metrics
  - function, 2718
- pgpro\_stats\_sdm\_stats\_updated
  - function, 2719
- pgpro\_stats\_sdm\_stats\_updated\_reset
  - function, 2719
- pgpro\_stats\_statements
  - function, 2718
- pgpro\_stats\_statements\_reset, 2718
- pgpro\_stats\_totals
  - function, 2718
- pgpro\_stats\_totals\_reset, 2718
- pgpro\_stats\_trace\_delete
  - function, 2727
- pgpro\_stats\_trace\_insert
  - function, 2727
- pgpro\_stats\_trace\_reset
  - function, 2727
- pgpro\_stats\_trace\_show
  - function, 2727
- pgpro\_stats\_trace\_update
  - function, 2727
- pgpro\_stats\_vacuum\_database
  - function, 2719
- pgpro\_stats\_vacuum\_indexes
  - function, 2719
- pgpro\_stats\_vacuum\_tables
  - function, 2719
- pgpro\_stats\_wal\_sender\_crc\_errors
  - function, 2719
- pgpro\_stat\_get\_wal\_activity, 834
- pgpro\_stat\_wal\_activity, 793, 810
- pgpro\_tune, 2098
- pgpro\_version, 370
- pgpro\_version configuration parameter, 664
- pgp\_armor\_headers, 2674
- pgp\_key\_id, 2673
- pgp\_pub\_decrypt, 2673
- pgp\_pub\_decrypt\_bytea, 2673
- pgp\_pub\_encrypt, 2673
- pgp\_pub\_encrypt\_bytea, 2673
- pgp\_sym\_decrypt, 2673
- pgp\_sym\_decrypt\_bytea, 2673
- pgp\_sym\_encrypt, 2673
- pgp\_sym\_encrypt\_bytea, 2673
- PGREQUIREAUTH, 978
- PGREQUIREPEER, 979
- PGREQUIRESSL, 979
- PGresult, 937
- PGREUSEPASS, 979
- PGSERVICE, 978
- PGSERVICEFILE, 979
- PGSSLCERT, 979
- PGSSLCERTMODE, 979
- PGSSLCOMPRESSION, 979
- PGSSLCRL, 979
- PGSSLCRLDIR, 979
- PGSSLKEY, 979
- PGSSLMAXPROTOCOLVERSION, 979
- PGSSLMINPROTOCOLVERSION, 979
- PGSSLMODE, 979
- PGSSLNEGOTIATION, 978
- PGSSLROOTCERT, 979
- PGSSLSNI, 979
- PGSYSCONFDIR, 979
- PGTARGETSESSIONATTRS, 979
- PGTZ, 979
- PGUSER, 978
- pgxs, 1229
- pg\_advisory\_lock, 387
- pg\_advisory\_lock\_shared, 387
- pg\_advisory\_unlock, 387
- pg\_advisory\_unlock\_all, 387
- pg\_advisory\_unlock\_shared, 388
- pg\_advisory\_xact\_lock, 388
- pg\_advisory\_xact\_lock\_shared, 388
- pg\_aggregate, 2208
- pg\_am, 2209
- pg\_amcheck, 1949
- pg\_amop, 2210
- pg\_amproc, 2211
- pg\_archivecleanup, 2085
- pg\_attrdef, 2211
- pg\_attribute, 2211
- pg\_authid, 2213
- pg\_auth\_members, 2215
- pg\_autoprepared\_statements, 2271

pg\_available\_extensions, 2261  
pg\_available\_extension\_versions, 2261  
pg\_available\_wal\_summaries, 370  
pg\_backend\_get\_config\_value, 371  
pg\_backend\_load\_library, 372  
pg\_backend\_memory\_contexts, 2262  
pg\_backend\_pid, 352  
pg\_backend\_set\_config, 371  
pg\_backup\_start, 375  
pg\_backup\_stop, 375  
pg\_basetype, 358  
pg\_blocking\_pids, 353  
pg\_cancel\_backend, 373  
pg\_cast, 2215  
pg\_char\_to\_encoding, 358  
pg\_checksums, 2087  
pg\_class, 2216  
pg\_client\_encoding, 213  
pg\_collation, 2218  
pg\_collation\_actual\_version, 383  
pg\_collation\_is\_visible, 357  
PG\_COLOR, 2974  
PG\_COLORS, 2974  
pg\_column\_compression, 382  
pg\_column\_size, 382  
pg\_column\_toast\_chunk\_id, 382  
pg\_combinebackup, 1977  
pg\_config, 1980, 2263  
    with ecpg, 1060  
    with libpq, 986  
    with user-defined C functions, 1181  
pg\_conf\_load\_time, 353  
pg\_constraint, 2219  
pg\_control\_checkpoint, 368  
pg\_control\_init, 368  
pg\_control\_recovery, 368  
pg\_control\_system, 368  
pg\_conversion, 2221  
pg\_conversion\_is\_visible, 357  
pg\_copy\_logical\_replication\_slot, 379  
pg\_copy\_physical\_replication\_slot, 379  
pg\_createsubscriber, 2089  
pg\_create\_logical\_replication\_slot, 379  
pg\_create\_physical\_replication\_slot, 379  
pg\_create\_restore\_point, 374  
pg\_ctl, 559, 561, 2093  
pg\_current\_logfile, 353  
pg\_current\_snapshot, 366  
pg\_current\_wal\_flush\_lsn, 374  
pg\_current\_wal\_insert\_lsn, 374  
pg\_current\_wal\_lsn, 374  
pg\_current\_xact\_id, 366  
pg\_current\_xact\_id\_if\_assigned, 366  
pg\_cursors, 2263  
pg\_database, 724, 2222  
pg\_database\_collation\_actual\_version, 383  
pg\_database\_size, 382  
pg\_db\_role\_setting, 2223  
pg\_ddl\_command, 196  
pg\_default\_acl, 2223  
pg\_depend, 2224  
pg\_describe\_object, 364  
pg\_description, 2226  
pg\_diagdump, 391  
pg\_drop\_replication\_slot, 379  
pg\_dump, 1983  
pg\_dumpall, 1997  
pg\_encoding\_to\_char, 358  
pg\_enum, 2226  
pg\_event\_trigger, 2227  
pg\_event\_trigger\_ddl\_commands, 393  
pg\_event\_trigger\_dropped\_objects, 393  
pg\_event\_trigger\_table\_rewrite\_oid, 395  
pg\_event\_trigger\_table\_rewrite\_reason, 395  
pg\_export\_snapshot, 378  
pg\_extension, 2227  
pg\_extension\_config\_dump, 1225  
pg\_filenode\_relation, 383  
pg\_file\_settings, 2264  
pg\_foreign\_data\_wrapper, 2228  
pg\_foreign\_server, 2228  
pg\_foreign\_table, 2229  
pg\_function\_is\_visible, 357  
pg\_get\_catalog\_foreign\_keys, 358  
pg\_get\_constraintdef, 359  
pg\_get\_expr, 359  
pg\_get\_functiondef, 359  
pg\_get\_function\_arguments, 359  
pg\_get\_function\_identity\_arguments, 359  
pg\_get\_function\_result, 359  
pg\_get\_indexdef, 359  
pg\_get\_keywords, 359  
pg\_get\_object\_address, 364  
pg\_get\_partkeydef, 359  
pg\_get\_ruledef, 359  
pg\_get\_serial\_sequence, 359  
pg\_get\_statisticsobjdef, 360  
pg\_get\_triggerdef, 360  
pg\_get\_userbyid, 360  
pg\_get\_viewdef, 360  
pg\_get\_wal\_replay\_pause\_state, 377  
pg\_get\_wal\_resource\_managers, 377  
pg\_get\_wal\_summarizer\_state, 370  
pg\_group, 2264  
pg\_has\_role, 356  
pg\_hba.conf, 689  
pg\_hba\_file\_rules, 2265  
pg\_ident.conf, 697  
pg\_identify\_object, 364  
pg\_identify\_object\_as\_address, 364  
pg\_ident\_file\_mappings, 2265  
pg\_import\_system\_collations, 383  
pg\_index, 2229  
pg\_indexam\_has\_property, 360  
pg\_indexes, 2266  
pg\_indexes\_size, 382

pg\_index\_column\_has\_property, 360  
pg\_index\_has\_property, 360  
pg\_inherits, 2231  
pg\_init\_privs, 2231  
pg\_input\_error\_info, 365  
pg\_input\_is\_valid, 365  
pg\_integrity\_check, 2903  
pg\_isready, 2004  
pg\_is\_in\_recovery, 376  
pg\_is\_other\_temp\_schema, 353  
pg\_is\_wal\_replay\_paused, 377  
pg\_jit\_available, 353  
pg\_language, 2232  
pg\_largeobject, 2232  
pg\_largeobject\_metadata, 2233  
pg\_last\_committed\_xact, 368  
pg\_last\_wal\_receive\_lsn, 376  
pg\_last\_wal\_replay\_lsn, 377  
pg\_last\_xact\_replay\_timestamp, 377  
pg\_listening\_channels, 353  
pg\_locks, 2266  
pg\_logical\_emit\_message, 381  
pg\_logical\_slot\_get\_binary\_changes, 380  
pg\_logical\_slot\_get\_changes, 379  
pg\_logical\_slot\_peek\_binary\_changes, 380  
pg\_logical\_slot\_peek\_changes, 380  
pg\_log\_backend\_memory\_contexts, 373  
pg\_log\_standby\_snapshot, 378  
pg\_lsn, 195  
pg\_ls\_archive\_statusdir, 386  
pg\_ls\_dir, 385  
pg\_ls\_logdir, 385  
pg\_ls\_logicalmapdir, 386  
pg\_ls\_logicalsnapdir, 386  
pg\_ls\_replslotdir, 386  
pg\_ls\_tmpdir, 386  
pg\_ls\_waldir, 385  
pg\_matviews, 2269  
pg\_mcv\_list\_items, 395  
pg\_my\_temp\_schema, 353  
pg\_namespace, 2233  
pg\_notification\_queue\_usage, 353  
pg\_notify, 1849  
pg\_opclass, 2234  
pg\_opclass\_is\_visible, 357  
pg\_operation\_log, 390  
pg\_operator, 2234  
pg\_operator\_is\_visible, 357  
pg\_opfamily, 2235  
pg\_opfamily\_is\_visible, 357  
pg\_options\_to\_table, 361  
pg\_parameter\_acl, 2236  
pg\_partitioned\_table, 2236  
pg\_partition\_ancestors, 384  
pg\_partition\_root, 384  
pg\_partition\_tree, 384  
pg\_policies, 2269  
pg\_policy, 2237  
pg\_postmaster\_start\_time, 353  
pg\_prepared\_statements, 2270  
pg\_prepared\_xacts, 2271  
pg\_proaudit.log\_catalog\_access configuration parameter, 2686  
pg\_proaudit.log\_command\_text configuration parameter, 2686  
pg\_proaudit.log\_destination configuration parameter, 2686  
pg\_proaudit.log\_directory configuration parameter, 2686  
pg\_proaudit.log\_filename configuration parameter, 2686  
pg\_proaudit.log\_rotation\_age configuration parameter, 2687  
pg\_proaudit.log\_rotation\_size configuration parameter, 2687  
pg\_proaudit.log\_truncate\_on\_rotation configuration parameter, 2687  
pg\_proaudit.max\_rules\_count configuration parameter, 2687  
pg\_probackup, 2837  
pg\_proc, 2237  
pg\_profile, 2240  
pg\_promote, 377  
pg\_publication, 2240  
pg\_publication\_namespace, 2241  
pg\_publication\_rel, 2241  
pg\_publication\_tables, 2272  
pg\_range, 2242  
pg\_read\_binary\_file, 387  
pg\_read\_file, 386  
pg\_receivewal, 2006  
pg\_receivexlog, 2975 (see [pg\\_receivewal](#))  
pg\_recvlogical, 2010  
pg\_relation\_filenode, 383  
pg\_relation\_filepath, 383  
pg\_relation\_size, 382  
pg\_reload\_conf, 373  
pg\_replication\_origin, 2242  
pg\_replication\_origin\_advance, 381  
pg\_replication\_origin\_create, 380  
pg\_replication\_origin\_drop, 380  
pg\_replication\_origin\_oid, 380  
pg\_replication\_origin\_progress, 381  
pg\_replication\_origin\_session\_is\_setup, 380  
pg\_replication\_origin\_session\_progress, 380  
pg\_replication\_origin\_session\_reset, 380  
pg\_replication\_origin\_session\_setup, 380  
pg\_replication\_origin\_status, 2272  
pg\_replication\_origin\_xact\_reset, 381  
pg\_replication\_origin\_xact\_setup, 381  
pg\_replication\_slots, 2273  
pg\_replication\_slot\_advance, 380  
pg\_resetwal, 2102  
pg\_resetxlog, 2975 (see [pg\\_resetwal](#))  
pg\_restore, 2014  
pg\_rewind, 2106

- pg\_rewrite, 2243
- pg\_roles, 2274
- pg\_role\_password, 2243
- pg\_rotate\_logfile, 373
- pg\_rules, 2276
- pg\_safe\_snapshot\_blocking\_pids, 353
- pg\_seclabel, 2243
- pg\_seclabels, 2276
- pg\_sequence, 2244
- pg\_sequences, 2276
- pg\_service.conf, 980
- pg\_settings, 2277
- pg\_settings\_get\_flags, 361
- pg\_shadow, 2279
- pg\_shdepend, 2244
- pg\_shdescription, 2246
- pg\_shmem\_allocations, 2280
- pg\_shseclabel, 2246
- pg\_size\_bytes, 382
- pg\_size\_pretty, 382
- pg\_sleep, 263
- pg\_sleep\_for, 263
- pg\_sleep\_until, 263
- pg\_snapshot\_any, 391
- pg\_snapshot\_xip, 366
- pg\_snapshot\_xmax, 366
- pg\_snapshot\_xmin, 366
- pg\_split\_walfile\_name, 375
- pg\_statio\_all\_indexes, 796, 826
- pg\_statio\_all\_sequences, 796, 827
- pg\_statio\_all\_tables, 796, 826
- pg\_statio\_sys\_indexes, 796
- pg\_statio\_sys\_sequences, 796
- pg\_statio\_sys\_tables, 796
- pg\_statio\_user\_indexes, 796
- pg\_statio\_user\_sequences, 796
- pg\_statio\_user\_tables, 796
- pg\_statistic, 490, 2246
- pg\_statistics\_obj\_is\_visible, 357
- pg\_statistic\_ext, 491, 2248
- pg\_statistic\_ext\_data, 491, 2248
- pg\_stats, 490, 2280
- pg\_stats\_ext, 2282
- pg\_stats\_ext\_exprs, 2283
- pg\_stats\_vacuum\_database, 2285
- pg\_stats\_vacuum\_indexes, 2286
- pg\_stats\_vacuum\_tables, 2288
- pg\_stat\_activity, 793, 796
- pg\_stat\_all\_indexes, 795, 825
- pg\_stat\_all\_tables, 795, 823
- pg\_stat\_archiver, 794, 817
- pg\_stat\_bgwriter, 794, 819
- pg\_stat\_checkpoint, 794, 820
- pg\_stat\_clear\_snapshot, 832
- pg\_stat\_database, 794, 821
- pg\_stat\_database\_conflicts, 795, 823
- pg\_stat\_file, 387
- pg\_stat\_get\_activity, 832
- pg\_stat\_get\_backend\_activity, 834
- pg\_stat\_get\_backend\_activity\_start, 834
- pg\_stat\_get\_backend\_client\_addr, 834
- pg\_stat\_get\_backend\_client\_port, 834
- pg\_stat\_get\_backend\_dbid, 834
- pg\_stat\_get\_backend\_idset, 834
- pg\_stat\_get\_backend\_pid, 834
- pg\_stat\_get\_backend\_start, 834
- pg\_stat\_get\_backend\_subxact, 834
- pg\_stat\_get\_backend\_userid, 834
- pg\_stat\_get\_backend\_wait\_event, 834
- pg\_stat\_get\_backend\_wait\_event\_type, 834
- pg\_stat\_get\_backend\_xact\_start, 834
- pg\_stat\_get\_snapshot\_timestamp, 832
- pg\_stat\_get\_xact\_blocks\_fetched, 832
- pg\_stat\_get\_xact\_blocks\_hit, 832
- pg\_stat\_gssapi, 794, 816
- pg\_stat\_io, 795, 817
- pg\_stat\_progress\_analyze, 794, 835
- pg\_stat\_progress\_basebackup, 794, 841
- pg\_stat\_progress\_cluster, 794, 836
- pg\_stat\_progress\_copy, 794, 837
- pg\_stat\_progress\_create\_index, 794, 838
- pg\_stat\_progress\_vacuum, 794, 840
- pg\_stat\_recovery\_prefetch, 794, 814
- pg\_stat\_replication, 793, 810
- pg\_stat\_replication\_slots, 795, 812
- pg\_stat\_reset, 832
- pg\_stat\_reset\_replication\_slot, 833
- pg\_stat\_reset\_shared, 832
- pg\_stat\_reset\_single\_function\_counters, 833
- pg\_stat\_reset\_single\_table\_counters, 832
- pg\_stat\_reset\_slru, 833
- pg\_stat\_reset\_subscription\_stats, 833
- pg\_stat\_slru, 795, 828
- pg\_stat\_ssl, 794, 816
- pg\_stat\_statements, 2740
  - function, 2744
- pg\_stat\_statements.max configuration parameter, 2745
- pg\_stat\_statements.save configuration parameter, 2745
- pg\_stat\_statements.track configuration parameter, 2745
- pg\_stat\_statements.track\_planning configuration parameter, 2745
- pg\_stat\_statements.track\_utility configuration parameter, 2745
- pg\_stat\_statements\_info, 2744
- pg\_stat\_statements\_reset, 2744
- pg\_stat\_subscription, 794, 815
- pg\_stat\_subscription\_stats, 795, 816
- pg\_stat\_sys\_indexes, 795
- pg\_stat\_sys\_tables, 795
- pg\_stat\_user\_functions, 795, 827
- pg\_stat\_user\_indexes, 795
- pg\_stat\_user\_tables, 795
- pg\_stat\_wal, 795, 820

- pg\_stat\_wal\_receiver, 794, 813
- pg\_stat\_xact\_all\_tables, 795
- pg\_stat\_xact\_sys\_tables, 795
- pg\_stat\_xact\_user\_functions, 795
- pg\_stat\_xact\_user\_tables, 795
- pg\_subscription, 2249
- pg\_subscription\_rel, 2250
- pg\_switch\_wal, 375
- pg\_sync\_replication\_slots, 381
- pg\_tables, 2289
- pg\_tablespace, 2251
- pg\_tablespace\_databases, 361
- pg\_tablespace\_location, 361
- pg\_tablespace\_size, 382
- pg\_table\_is\_visible, 357
- pg\_table\_size, 382
- pg\_temp, 646
  - securing functions, 1626
- pg\_temp\_relation\_size, 382
- pg\_terminate\_backend, 373
- pg\_test\_fsync, 2112
- pg\_test\_timing, 2113
- pg\_timezone\_abbrevs, 2290
- pg\_timezone\_names, 2290
- pg\_total\_relation\_size, 383
- pg\_transform, 2251
- pg\_trigger, 2252
- pg\_trigger\_depth, 354
- pg\_try\_advisory\_lock, 388
- pg\_try\_advisory\_lock\_shared, 388
- pg\_try\_advisory\_xact\_lock, 388
- pg\_try\_advisory\_xact\_lock\_shared, 388
- pg\_ts\_config, 2253
- pg\_ts\_config\_is\_visible, 357
- pg\_ts\_config\_map, 2254
- pg\_ts\_dict, 2254
- pg\_ts\_dict\_is\_visible, 358
- pg\_ts\_parser, 2254
- pg\_ts\_parser\_is\_visible, 358
- pg\_ts\_template, 2255
- pg\_ts\_template\_is\_visible, 358
- pg\_type, 2255
- pg\_typeof, 361
- pg\_type\_is\_visible, 358
- pg\_user, 2290
- pg\_user\_mapping, 2259
- pg\_user\_mappings, 2291
- pg\_verifybackup, 2025
- pg\_views, 2292
- pg\_visible\_in\_snapshot, 366
- pg\_wait\_events, 2292
- pg\_wait\_sampling.history\_period variable, 2749
- pg\_wait\_sampling.history\_size variable, 2749
- pg\_wait\_sampling.profile\_period variable, 2750
- pg\_wait\_sampling.profile\_pid variable, 2750
- pg\_wait\_sampling.profile\_queries variable, 2750
- pg\_wait\_sampling.sample\_cpu variable, 2750
- pg\_wait\_sampling.get\_current, 2751
- pg\_wait\_sampling.reset\_profile, 2751
- pg\_waldump, 2116
- pg\_walfile\_name, 375
- pg\_walfile\_name\_offset, 375
- pg\_walsummary, 2120
- pg\_wal\_lsn\_diff, 375
- pg\_wal\_replay\_pause, 377
- pg\_wal\_replay\_resume, 377
- pg\_wal\_summary\_contents, 370
- pg\_xact\_commit\_timestamp, 368
- pg\_xact\_commit\_timestamp\_origin, 368
- pg\_xact\_status, 366
- pg\_xlogdump, 2975 (see [pg\\_waldump](#))
- phantom read, 462
- phraseto\_tsquery, 275, 433
- pi, 206
- PIC, 1181
- PID
  - determining PID of server process
    - in libpq, 931
- pipelining
  - in libpq, 952
  - protocol specification, 2310
- PITR, 761
- PITR standby, 770
- pkg-config
  - with ecpg, 1060
  - with libpq, 986
- PL/Perl, 1350
- PL/PerlU, 1360
- PL/pgSQL, 1278
- PL/Python, 1365
- PL/SQL (Oracle)
  - porting to PL/pgSQL, 1331
- PL/Tcl, 1340
- plainto\_tsquery, 275, 433
- planner\_upper\_limit\_estimation configuration parameter, 622
- plan\_cache\_lru\_memsize configuration parameter, 596
- plan\_cache\_lru\_size configuration parameter, 597
- plan\_cache\_mode configuration parameter, 627
- plperl.on\_init configuration parameter, 1363
- plperl.on\_plperl\_init configuration parameter, 1363
- plperl.on\_plperl\_init configuration parameter, 1363
- plperl.use\_strict configuration parameter, 1363
- plpgsql.check\_asserts configuration parameter, 1316
- plpgsql.variable\_conflict configuration parameter, 1326
- pltcl.start\_proc configuration parameter, 1349
- pltclu.start\_proc configuration parameter, 1349
- point, 156, 270
- point-in-time recovery, 761
- policy, 75
- polygon, 157, 270
- polymorphic function, 1154

- polymorphic type, 1154
- popcount (see [bit\\_count](#))
- popen, 269
- port, 920
- port configuration parameter, 675
- portal
  - DECLARE, 1755
  - in PL/pgSQL, 1309
- position, 211, 220, 224
- POSITION\_REGEX, 240
- POSTGRES, xxv
- postgres, 11, 560, 2121
- postgres user, 558
- Postgres95, xxvi
- postgresql.auto.conf, 582
- postgresql.conf, 581
- postgres\_fdw, 2752
- postgres\_fdw.additional\_ordered\_paths configuration parameter, 675
- postgres\_fdw.application\_name configuration parameter, 2760
- postgres\_fdw.enforce\_foreign\_join configuration parameter, 675
- postgres\_fdw.foreign\_explain configuration parameter, 675
- postgres\_fdw.remote\_plan\_cache configuration parameter, 681
- postgres\_fdw.use\_twophase configuration parameter, 675
- post\_auth\_delay configuration parameter, 666
- power, 206
- PQbackendPID, 931
- PQbinaryTuples, 943
  - with COPY, 963
- PQcancel, 960
- PQcancelBlocking, 957
- PQcancelCreate, 957
- PQcancelErrorMessage, 959
- PQcancelFinish, 960
- PQcancelPoll, 958
- PQcancelReset, 960
- PQcancelSocket, 959
- PQcancelStart, 958
- PQcancelStatus, 959
- PQchangePassword, 969
- PQclear, 941
- PQclientEncoding, 966
- PQclosePortal, 937
- PQclosePrepared, 937
- PQcmdStatus, 945
- PQcmdTuples, 945
- PQconndefaults, 914
- PQconnectdb, 911
- PQconnectdbParams, 910
- PQconnectionNeedsPassword, 932
- PQconnectionUsedGSSAPI, 932
- PQconnectionUsedPassword, 932
- PQconnectPoll, 912
- PQconnectStart, 912
- PQconnectStartParams, 912
- PQconninfo, 915
- PQconninfoFree, 968
- PQconninfoParse, 915
- PQconsumeInput, 950
- PQcopyResult, 970
- PQdb, 928
- PQdescribePortal, 937
- PQdescribePrepared, 936
- PQencryptPassword, 969
- PQencryptPasswordConn, 968
- PQendcopy, 966
- PQenterPipelineMode, 954
- PQerrorMessage, 931
- PQescapeBytea, 947
- PQescapeByteaConn, 947
- PQescapeIdentifier, 946
- PQescapeLiteral, 945
- PQescapeString, 947
- PQescapeStringConn, 946
- PQexec, 934
- PQexecParams, 934
- PQexecPrepared, 936
- PQexitPipelineMode, 954
- PQfformat, 942
  - with COPY, 963
- PQfinish, 915
- PQfireResultCreateEvents, 970
- PQflush, 952
- PQfmod, 943
- PQfn, 961
- PQfname, 942
- PQfnumber, 942
- PQfreeCancel, 960
- PQfreemem, 968
- PQfsize, 943
- PQftable, 942
- PQftablecol, 942
- PQftype, 943
- PQgetCancel, 960
- PQgetCopyData, 964
- PQgetcurrentTimeUSec, 971
- PQgetisnull, 944
- PQgetlength, 944
- PQgetline, 965
- PQgetlineAsync, 965
- PQgetResult, 950
- PQgetssl, 933
- PQgetSSLKeyPassHook\_OpenSSL, 917
- PQgetvalue, 943
- PQhost, 929
- PQhostaddr, 929
- PQinitOpenSSL, 985
- PQinitSSL, 985
- PQinstanceData, 975
- PQisBusy, 951
- PQisnonblocking, 951



- PQisthreadsafe, 985
- PQlibVersion, 971
  - (see also [PQserverVersion](#))
- PQmakeEmptyPGresult, 969
- PQnfields, 941
  - with COPY, 963
- PQnotifies, 962
- PQnparams, 944
- PQntuples, 941
- PQoidStatus, 945
- PQoidValue, 945
- PQoptions, 930
- PQparameterStatus, 930
- PQparamtype, 944
- PQpass, 928
- PQping, 916
- PQpingParams, 916
- PQpipelineStatus, 954
- PQpipelineSync, 955
- PQport, 929
- PQprepare, 935
- PQprint, 944
- PQprotocolVersion, 931
- PQputCopyData, 963
- PQputCopyEnd, 964
- PQputline, 965
- PQputnbytes, 966
- PQregisterEventProc, 975
- PQrequestCancel, 961
- PQreset, 915
- PQresetPoll, 915
- PQresetStart, 915
- PQresStatus, 938
- PQresultAlloc, 970
- PQresultErrorField, 939
- PQresultErrorMessage, 939
- PQresultInstanceData, 975
- PQresultMemorySize, 971
- PQresultSetInstanceData, 975
- PQresultStatus, 937
- PQresultVerboseErrorMessage, 939
- PQsendClosePortal, 950
- PQsendClosePrepared, 950
- PQsendDescribePortal, 949
- PQsendDescribePrepared, 949
- PQsendFlushRequest, 955
- PQsendPipelineSync, 955
- PQsendPrepare, 949
- PQsendQuery, 948
- PQsendQueryParams, 948
- PQsendQueryPrepared, 949
- PQserverVersion, 931
- PQsetChunkedRowsMode, 956
- PQsetClientEncoding, 967
- PQsetdb, 911
- PQsetdbLogin, 911
- PQsetErrorContextVisibility, 967
- PQsetErrorVerbosity, 967
- PQsetInstanceData, 975
- PQsetnonblocking, 951
- PQsetNoticeProcessor, 972
- PQsetNoticeReceiver, 972
- PQsetResultAttrs, 970
- PQsetSingleRowMode, 956
- PQsetSSLKeyPassHook\_OpenSSL, 916
- PQsetTraceFlags, 968
- PQsetvalue, 970
- PQsocket, 931
- PQsocketPoll, 914
- PQsslAttribute, 932
- PQsslAttributeNames, 933
- PQsslInUse, 932
- PQsslStruct, 933
- PQstatus, 930
- PQtrace, 967
- PQtransactionStatus, 930
- PQtty, 929
- PQunescapeBytea, 948
- PQuntrace, 968
- PQuser, 928
- predicate locking, 466
- PREPARE, 1850
- PREPARE TRANSACTION, 1853
- prepared statements
  - creating, 1850
  - executing, 1809
  - removing, 1754
  - showing the query plan, 1810
- preparing a query
  - in PL/pgSQL, 1327
  - in PL/Python, 1373
  - in PL/Tcl, 1343
- pre\_auth\_delay configuration parameter, 666
- primary key, 64
- primary\_conninfo configuration parameter, 615
- primary\_slot\_name configuration parameter, 616
- privilege, 71
  - default, 73
  - querying, 354
  - with rules, 1271
  - for schemas, 85
  - with views, 1271
- procedural language, 1276
  - externally maintained, 2938
  - handler for, 2354
- procedure
  - user-defined, 1157
- procedures
  - output parameter, 1163
- protocol
  - frontend-backend, 2301
- ps
  - to monitor activity, 791
- psql, 13, 2028
- ptrack.map\_size configuration parameter, 2763
- ptrack\_get\_change\_stat, 2763



ptrack\_get\_pagemapset, 2763  
ptrack\_version, 2763  
Python, 1365

## Q

qualified name, 82  
query, 18, 107  
query cancellation, 957  
query plan, 477  
query tree, 1252  
querytree, 276, 439  
quotation marks  
    and identifiers, 34  
    escaping, 35  
quote\_all\_identifiers configuration parameter, 658  
quote\_ident, 213  
    in PL/Perl, 1358  
    use in PL/pgSQL, 1291  
quote\_literal, 213  
    in PL/Perl, 1358  
    use in PL/pgSQL, 1292  
quote\_nullable, 214  
    in PL/Perl, 1358  
    use in PL/pgSQL, 1292

## R

radians, 206  
radius, 269  
RADIUS, 706  
RAISE  
    in PL/pgSQL, 1314  
random, 207, 207  
random\_normal, 207  
random\_page\_cost configuration parameter, 622  
range table, 1252  
range type, 188  
    exclude, 192  
    indexes on, 192  
range\_agg, 337  
range\_intersect\_agg, 337  
range\_merge, 333, 334  
rank, 342  
    hypothetical, 340  
read committed, 463  
read-only transaction, 648  
    setting, 1908  
    setting default, 647  
Readline  
    in psql, 2063  
READ\_REPLICATION\_SLOT, 2319  
real, 137  
REASSIGN OWNED, 1855  
rebalancing, 894  
record, 196  
recovery.conf, 2975, 2975  
recovery.signal, 609  
recovery\_end\_command configuration parameter, 610

recovery\_init\_sync\_method configuration parameter, 662  
recovery\_min\_apply\_delay configuration parameter, 618  
recovery\_prefetch configuration parameter, 608  
recovery\_target configuration parameter, 610  
recovery\_target\_action configuration parameter, 611  
recovery\_target\_inclusive configuration parameter, 611  
recovery\_target\_lsn configuration parameter, 611  
recovery\_target\_name configuration parameter, 610  
recovery\_target\_time configuration parameter, 610  
recovery\_target\_timeline configuration parameter, 611  
recovery\_target\_xid configuration parameter, 611  
rectangle, 156  
RECURSIVE  
    in common table expressions, 126  
    in views, 1749  
recursive\_worktable\_factor configuration parameter, 627  
referential integrity, 25, 65  
REFRESH MATERIALIZED VIEW, 1856  
regclass, 193  
regcollation, 193  
regconfig, 193  
regdictionary, 193  
regexp\_count, 214, 227  
regexp\_instr, 214, 227  
regexp\_like, 214, 227  
regexp\_match, 214, 227  
regexp\_matches, 214, 227  
regexp\_replace, 214, 227  
regexp\_split\_to\_array, 215, 227  
regexp\_split\_to\_table, 215, 227  
regexp\_substr, 215, 227  
regnamespace, 193  
regoper, 193  
regoperator, 193  
regproc, 193  
regprocedure, 193  
regprofile, 193  
regression intercept, 338  
regression slope, 338  
regrole, 193  
regr\_avgx, 338  
regr\_avgy, 338  
regr\_count, 338  
regr\_intercept, 338  
regr\_r2, 338  
regr\_slope, 338  
regr\_sxx, 339  
regr\_sxy, 339  
regr\_syy, 339  
regtype, 193  
regular expression, 226, 227

- (see also [pattern matching](#))
  - regular expressions
    - and locales, 728
  - reindex, 759
  - REINDEX, 1858
  - reindexdb, 2070
  - relation, 16
  - relational database, 16
  - RELEASE SAVEPOINT, 1863
  - REMOTE EXPLAIN parameter, 504
  - remove\_temp\_files\_after\_crash configuration parameter, 669
  - rename\_server, 2772
  - repeat, 215
  - repeatable read, 465
  - replace, 215
  - replication, 770
  - Replication Origins, 1462
  - Replication Progress Tracking, 1462
  - replication slot
    - logical replication, 1450
    - streaming replication, 777
  - reporting errors
    - in PL/pgSQL, 1314
  - reserved\_connections configuration parameter, 586
  - RESET, 1865
  - restartpoint, 850
  - restart\_after\_crash configuration parameter, 661
  - restore\_command configuration parameter, 609
  - RESTRICT
    - with DROP, 102
    - foreign key action, 66
  - restrict\_nonsystem\_relation\_kind configuration parameter, 652
  - retryable error, 475
  - RETURN NEXT
    - in PL/pgSQL, 1295
  - RETURN QUERY
    - in PL/pgSQL, 1295
  - RETURNING, 106
  - RETURNING INTO
    - in PL/pgSQL, 1289
  - reverse, 215
  - REVOKE, 71, 1866
  - right, 215
  - right join, 109
  - role, 710, 714
    - applicable, 1105
    - enabled, 1122
    - membership in, 712
    - privilege to bypass, 712
    - privilege to create, 711
    - privilege to inherit, 712
    - privilege to initiate replication, 711
    - privilege to limit connection, 712
  - ROLLBACK, 1870
  - rollback
    - psql, 2059
  - ROLLBACK PREPARED, 1871
  - ROLLBACK TO SAVEPOINT, 1872
  - ROLLUP, 118
  - round, 206
  - routine, 1157
  - routine maintenance, 750
  - row, 16, 56
  - ROW, 51
  - row estimation
    - multivariate, 2481
    - planner, 2477
  - row type, 182
    - constructor, 51
  - row-level security, 75
  - row-wise comparison, 346
  - row\_number, 341
  - row\_security configuration parameter, 646
  - row\_security\_active, 356
  - row\_to\_json, 296
  - rpad, 211
  - rtrim, 211, 220
  - rule, 1252
    - and materialized views, 1260
    - and views, 1253
    - for DELETE, 1262
    - for INSERT, 1262
    - for SELECT, 1254
    - compared with triggers, 1274
    - for UPDATE, 1262
- ## S
- SAVEPOINT, 1874
  - savepoints
    - defining, 1874
    - releasing, 1863
    - rolling back, 1872
  - scalar (see [expression](#))
  - scale, 206
  - scaling, 894
  - schema, 81
    - creating, 82
    - current, 84, 352
    - public, 83
    - removing, 82
    - vault, 83
  - SCRAM, 699
  - scram\_iterations configuration parameter, 588
  - search path, 83
    - current, 352
    - object visibility, 357
  - search\_path configuration parameter, 84, 645
    - use in securing functions, 1626
  - secure communication, 894
  - SECURITY LABEL, 1876
  - segment\_size configuration parameter, 664
  - SELECT, 18, 107, 1878
    - determination of result type, 409

- select list, 121
- SELECT INTO, 1898
  - in PL/pgSQL, 1289
- self\_join\_search\_limit configuration parameter, 622
- semaphores, 564
- send\_abort\_for\_crash configuration parameter, 669
- send\_abort\_for\_kill configuration parameter, 669
- sequence, 322
  - and serial type, 138
- sequential scan, 622
- seq\_page\_cost configuration parameter, 622
- seq\_scan\_startup\_cost\_first\_row configuration parameter, 627
- serial, 138
- serial2, 138
- serial4, 138
- serial8, 138
- serializable, 466
- Serializable Snapshot Isolation, 462
- serializable\_buffers configuration parameter, 595
- serialization anomaly, 462, 466
- serialization failure, 475
- server log, 627
  - log file maintenance, 759
- Server Name Indication, 926
- server spoofing, 574
- server\_encoding configuration parameter, 664
- server\_version configuration parameter, 664
- server\_version\_num configuration parameter, 664
- session\_cpu\_weight configuration parameter, 601
- session\_ioread\_weight configuration parameter, 601
- session\_iowrite\_weight configuration parameter, 602
- session\_preload\_libraries configuration parameter, 654
- session\_replication\_role configuration parameter, 648
- session\_user, 354
- SET, 371, 1900
- SET CONSTRAINTS, 1903
- set difference, 122
- set intersection, 122
- set operation, 122
- set returning functions
  - functions, 348
- SET ROLE, 1904
- SET SESSION AUTHORIZATION, 1906
- SET TRANSACTION, 1908
- set union, 122
- SET XML OPTION, 651
- setseed, 207
- setval, 322
- setweight, 276, 438
  - setweight for specific lexeme(s), 276
- set\_bit, 221, 224
- set\_byte, 221
- set\_config, 371
- set\_masklen, 273
- set\_server\_connstr, 2772
- set\_server\_db\_exclude, 2772
- set\_server\_description, 2772
- set\_server\_max\_sample\_age, 2772
- set\_server\_setting, 2773
- set\_server\_size\_sampling, 2775
- set\_server\_subsampling, 2772
- sha224, 221
- sha256, 221
- sha384, 221
- sha512, 221
- shardman.advisory\_xact\_lock, 388
- shardman.advisory\_xact\_lock\_shared, 388
- shardman.am\_coordinator() , 398
- shardman.attach\_subpart , 397
- shardman.broadcast\_all\_sql , 396
- shardman.broadcast\_ddl configuration parameter, 675
- shardman.broadcast\_query , 396
- shardman.broadcast\_sql , 396
- shardman.context\_log configuration parameter, 637
- shardman.create\_subpart , 397
- shardman.database configuration parameter, 680
- shardman.detach\_subpart , 397
- shardman.drop\_subpart , 397
- shardman.enable\_limit\_pushdown configuration parameter, 675
- shardman.get\_partition\_for\_value , 396
- shardman.global\_analyze , 397
- shardman.gt\_batch\_size configuration parameter, 681
- shardman.monitor\_deadlock\_interval configuration parameter, 681
- shardman.monitor\_dxact\_interval configuration parameter, 680
- shardman.monitor\_dxact\_timeout configuration parameter, 681
- shardman.monitor\_interval configuration parameter, 680
- shardman.monitor\_trim\_csnxid\_map\_interval configuration parameter, 681
- shardman.num\_parts configuration parameter, 676
- shardman.oldest\_csn, 830
- shardman.pg\_indoubt\_xacts, 829
- shardman.pg\_stat\_csn, 829
- shardman.pg\_stat\_foreign\_stat\_bytes, 831
- shardman.pg\_stat\_monitor, 830
- shardman.pg\_stat\_netusage, 831
- shardman.pg\_stat\_xact\_time, 830
- shardman.plan\_cache\_mem configuration parameter, 681
- shardman.query\_engine\_mode configuration parameter, 676
- shardman.rgid configuration parameter, 676
- shardman.silkroad\_sched\_priority configuration parameter, 679

- shardman.silkworm\_fetch\_size configuration parameter, 677
- shardman.silkworm\_sched\_priority configuration parameter, 679
- shardman.silk\_backends, 2494
- shardman.silk\_connects, 2493
- shardman.silk\_flow\_control configuration parameter, 680
- shardman.silk\_hello\_timeout configuration parameter, 678
- shardman.silk\_listen\_ip configuration parameter, 677
- shardman.silk\_max\_message configuration parameter, 677
- shardman.silk\_num\_workers configuration parameter, 677
- shardman.silk\_pending\_jobs, 2497
- shardman.silk\_rbc\_snap() , 2500
- shardman.silk\_routes, 2492
- shardman.silk\_routing, 2495
- shardman.silk\_routing() , 2500
- shardman.silk\_scheduler\_mode configuration parameter, 678
- shardman.silk\_set\_affinity configuration parameter, 679
- shardman.silk\_state, 2499
- shardman.silk\_statinfo, 2497
- shardman.silk\_statinfo\_reset() , 2500
- shardman.silk\_stream\_work\_mem configuration parameter, 677
- shardman.silk\_tracelog configuration parameter, 680
- shardman.silk\_tracelog\_category configuration parameter, 680
- shardman.silk\_tracepoints configuration parameter, 677
- shardman.silk\_track\_time configuration parameter, 680
- shardman.silk\_unassigned\_job\_queue\_size configuration parameter, 677
- shardman.silk\_use\_ip configuration parameter, 676
- shardman.silk\_use\_port configuration parameter, 677
- shardman.sync\_cluster\_settings configuration parameter, 676
- shardman.sync\_cluster\_settings\_blacklist configuration parameter, 676
- shardman.sync\_schema configuration parameter, 676
- shardman.trim\_csnxid\_map\_naptime configuration parameter, 681
- shardman.try\_advisory\_xact\_lock, 389
- shardman.try\_advisory\_xact\_lock\_shared, 389
- shardmanctl, 2128
- shardmand, 2194
- shared library, 1181
- shared memory, 564
- shared\_buffers configuration parameter, 591
- shared\_memory\_size configuration parameter, 664
- shared\_memory\_size\_in\_huge\_pages configuration parameter, 664
- shared\_memory\_type configuration parameter, 595
- shared\_preload\_libraries, 1191, 1192
- shared\_preload\_libraries configuration parameter, 655
- shobj\_description, 365
- SHOW, 371, 1911, 2317
- show\_baselines, 2778
- show\_samples, 2774
- show\_servers, 2773
- show\_server\_settings, 2773
- shutdown, 571
- SIGHUP, 582, 689, 697
- SIGINT, 572
- sign, 206
- signal
  - backend processes, 373
- significant digits, 652
- SIGQUIT, 572
- SIGTERM, 572
- SIMILAR TO, 226
- sin, 209
- sind, 209
- single-user mode, 2124
- sinh, 209
- skip\_temp\_rel\_lock
  - configuration parameter, 652
- sleep, 263
- sliced bread (see [TOAST](#))
- slope, 269
- SLRU, 828
- smallint, 135
- smallserial, 138
- Solaris
  - shared library, 1182
  - start script, 562
- SOME, 338, 343, 346
- sorting, 123
- SP-GiST (see [index](#))
- SPI, 1379
- spi\_commit
  - in PL/Perl, 1357
- SPI\_commit, 1438
- SPI\_commit\_and\_chain, 1438
- SPI\_connect, 1380
- SPI\_connect\_ext, 1380
- SPI\_copypuple, 1431
- spi\_cursor\_close
  - in PL/Perl, 1355
- SPI\_cursor\_close, 1413
- SPI\_cursor\_fetch, 1409
- SPI\_cursor\_find, 1408
- SPI\_cursor\_move, 1410
- SPI\_cursor\_open, 1403
- SPI\_cursor\_open\_with\_args, 1404
- SPI\_cursor\_open\_with\_paramlist, 1406

- SPI\_cursor\_parse\_open, 1407
- SPI\_exec, 1385
- SPI\_execp, 1402
- SPI\_execute, 1382
- SPI\_execute\_extended, 1386
- SPI\_execute\_plan, 1398
- SPI\_execute\_plan\_extended, 1399
- SPI\_execute\_plan\_with\_paramlist, 1401
- SPI\_execute\_with\_args, 1388
- spi\_exec\_prepared
  - in PL/Perl, 1356
- spi\_exec\_query
  - in PL/Perl, 1354
- spi\_fetchrow
  - in PL/Perl, 1355
- SPI\_finish, 1381
- SPI\_fname, 1419
- SPI\_fnumber, 1420
- spi\_freeplan
  - in PL/Perl, 1356
- SPI\_freeplan, 1437
- SPI\_freetuple, 1435
- SPI\_freetuptable, 1436
- SPI\_getargcount, 1395
- SPI\_getargtypeid, 1396
- SPI\_getbinval, 1422
- SPI\_getnspname, 1426
- SPI\_getrelname, 1425
- SPI\_gettype, 1423
- SPI\_gettypeid, 1424
- SPI\_getvalue, 1421
- SPI\_is\_cursor\_plan, 1397
- SPI\_keepplan, 1414
- SPI\_modifytuple, 1433
- SPI\_palloc, 1428
- SPI\_pfree, 1430
- spi\_prepare
  - in PL/Perl, 1356
- SPI\_prepare, 1390
- SPI\_prepare\_cursor, 1392
- SPI\_prepare\_extended, 1393
- SPI\_prepare\_params, 1394
- spi\_query
  - in PL/Perl, 1355
- spi\_query\_prepared
  - in PL/Perl, 1356
- SPI\_register\_relation, 1416
- SPI\_register\_trigger\_data, 1418
- SPI\_repallocc, 1429
- SPI\_result\_code\_string, 1427
- SPI\_returntuple, 1432
- spi\_rollback
  - in PL/Perl, 1357
- SPI\_rollback, 1439
- SPI\_rollback\_and\_chain, 1439
- SPI\_saveplan, 1415
- SPI\_scroll\_cursor\_fetch, 1411
- SPI\_scroll\_cursor\_move, 1412
- SPI\_start\_transaction, 1440
- SPI\_unregister\_relation, 1417
- split\_part, 215
- SQL/CLI, 2548
- SQL/Foundation, 2548
- SQL/Framework, 2548
- SQL/JRT, 2548
- SQL/JSON
  - functions and expressions, 292
- SQL/JSON path language, 304
- SQL/MDA, 2548
- SQL/MED, 2548
- SQL/OLB, 2548
- SQL/PGQ, 2548
- SQL/PSM, 2548
- SQL/Schemata, 2548
- SQL/XML, 2548
  - limits and conformance, 2569
- sqrt, 206
- ssh, 579
- SSI, 462
- SSL
  - in libpq, 933
  - with libpq, 924
  - TLS, 575, 981
- ssl configuration parameter, 589
- ssl\_ca\_file configuration parameter, 589
- ssl\_cert\_file configuration parameter, 589
- ssl\_ciphers configuration parameter, 590
- ssl\_crl\_dir configuration parameter, 589
- ssl\_crl\_file configuration parameter, 589
- ssl\_dh\_params\_file configuration parameter, 591
- ssl\_ecdh\_curve configuration parameter, 590
- ssl\_key\_file configuration parameter, 589
- ssl\_library configuration parameter, 664
- ssl\_max\_protocol\_version configuration parameter, 591
- ssl\_min\_protocol\_version configuration parameter, 590
- ssl\_passphrase\_command configuration parameter, 591
- ssl\_passphrase\_command\_supports\_reload configuration parameter, 591
- ssl\_prefer\_server\_ciphers configuration parameter, 590
- SSPI, 701
- STABLE, 1172
- standard deviation, 339
  - population, 339
  - sample, 339
- standard\_conforming\_strings configuration parameter, 658
- standby server, 770
- standby.signal, 609, 774, 775
  - for hot standby, 787
- standby\_mode (see [standby.signal](#))
- START TRANSACTION, 1913
- starts\_with, 215

- START\_REPLICATION, 2319
  - statement\_timeout configuration parameter, 648
  - statement\_timestamp, 254
  - statistics, 338, 792
    - of the planner, 489, 491, 752
  - stats\_fetch\_consistency configuration parameter, 642
  - stddev, 339
  - stddev\_pop, 339
  - stddev\_samp, 339
  - STONITH, 770
  - storage parameters, 1542, 1707, 1707
  - Streaming Replication, 770
  - string (see [character string](#))
  - strings
    - backslash quotes, 658
    - escape warning, 658
    - standard conforming, 658
  - string\_agg, 337
  - string\_to\_array, 215
  - string\_to\_table, 216
  - strip, 276, 438
  - strpos, 216
  - subplan, 483
    - hashed, 484
  - subquery, 22, 49, 112, 343
  - subscript, 42
  - substr, 216, 221
  - substring, 211, 220, 224, 226, 227
  - SUBSTRING\_REGEX, 240
  - subtransactions
    - in PL/Tcl, 1347
  - subtransaction\_buffers configuration parameter, 595
  - sum, 337
  - summarize\_wal configuration parameter, 612
  - superuser, 13, 711
  - superuser\_reserved\_connections configuration parameter, 586
  - support functions
    - in\_range, 2407
  - suppress\_redundant\_updates\_trigger, 392
  - synchronized\_standby\_slots configuration parameter, 614
  - synchronize\_seqscans configuration parameter, 659
  - synchronous commit, 847
  - Synchronous Replication, 770
  - synchronous\_commit configuration parameter, 603
  - synchronous\_standby\_gap configuration parameter, 615
  - synchronous\_standby\_names configuration parameter, 614
  - sync\_replication\_slots configuration parameter, 618
  - syntax
    - SQL, 33
  - syslog\_facility configuration parameter, 630
  - syslog\_ident configuration parameter, 630
  - syslog\_sequence\_numbers configuration parameter, 630
  - syslog\_split\_messages configuration parameter, 630
  - system catalog
    - schema, 85
  - systemd, 561
    - RemoveIPC, 567
  - system\_user, 354
- ## T
- table, 16, 56
    - creating, 56
    - inheritance, 86
    - modifying, 69
    - partitioning, 89
    - removing, 57
    - renaming, 71
  - Table Access Method, 2387
  - TABLE command, 1878
  - table expression, 107
  - table function, 112
    - XMLTABLE, 286
  - table sampling method, 2375
  - tableam
    - Table Access Method, 2387
  - tableoid, 68
  - TABLESAMPLE method, 2375
  - tablespace, 725
    - default, 646
    - temporary, 647
  - table\_am\_handler, 196
  - take\_sample, 2773
  - take\_sample\_subset, 2773, 2777
  - take\_subsample, 2776
  - tan, 209
  - tand, 209
  - tanh, 209
  - target list, 1253
  - Tcl, 1340
  - tcp\_keepalives\_count configuration parameter, 587
  - tcp\_keepalives\_idle configuration parameter, 587
  - tcp\_keepalives\_interval configuration parameter, 587
  - tcp\_user\_timeout configuration parameter, 588
  - template0, 723
  - template1, 723
  - temp\_buffers configuration parameter, 592
  - temp\_file\_limit configuration parameter, 596
  - temp\_tablespace configuration parameter, 647
  - temp\_table\_max\_size configuration parameter, 597
  - text, 140, 273
  - text search, 426
    - data types, 161
    - functions and operators, 161
    - indexes, 457
  - threads

- with libpq, 985
  - tid, 193
  - time, 144, 146
    - constants, 148
    - current, 262
    - output format, 148
      - (see also [formatting](#))
  - time span, 144
  - time with time zone, 144, 146
  - time without time zone, 144, 146
  - time zone, 149, 652
    - conversion, 260
    - input abbreviations, 2516
    - POSIX-style specification, 2518
  - time zone names, 652
  - timelines, 761
  - TIMELINE\_HISTORY, 2318
  - timeofday, 254
  - timeout
    - client authentication, 588
    - deadlock, 656
  - timestamp, 144, 147
  - timestamp with time zone, 144, 147
  - timestamp without time zone, 144, 147
  - timestampptz, 144
  - TimeZone configuration parameter, 652
  - timezone\_abbreviations configuration parameter, 652
  - to
    - storage parameter, 1633
  - TOAST, 2461
    - and user-defined types, 1204
    - per-column storage settings, 1533, 1698
    - per-type storage settings, 1558
    - versus large objects, 997
  - toast\_tuple\_target storage parameter, 1708
  - token, 33
  - to\_ascii, 216
  - to\_bin, 216
  - to\_char, 242
    - and locales, 728
  - to\_date, 242
  - to\_hex, 216
  - to\_json, 296
  - to\_jsonb, 296
  - to\_number, 242
  - to\_oct, 216
  - to\_regclass, 361
  - to\_regcollation, 361
  - to\_regnamespace, 361
  - to\_regoper, 361
  - to\_regoperator, 361
  - to\_regproc, 362
  - to\_regprocedure, 362
  - to\_regprofile, 362
  - to\_regrole, 362
  - to\_regtype, 362
  - to\_regtypemod, 362
- to\_timestamp, 242, 254
  - to\_tsquery, 276, 432
  - to\_tsvector, 276, 431
  - trace\_locks configuration parameter, 667
  - trace\_lock\_oidmin configuration parameter, 667
  - trace\_lock\_table configuration parameter, 667
  - trace\_lwlocks configuration parameter, 667
  - trace\_notify configuration parameter, 666
  - trace\_sort configuration parameter, 666
  - trace\_userlocks configuration parameter, 667
  - track\_activities configuration parameter, 641
  - track\_activity\_query\_size configuration parameter, 642
  - track\_commit\_timestamp configuration parameter, 613
  - track\_counts configuration parameter, 642
  - track\_fdw\_wait\_timing configuration parameter, 682
  - track\_functions configuration parameter, 642
  - track\_io\_timing configuration parameter, 642
  - track\_wal\_io\_timing configuration parameter, 642
  - track\_xact\_time configuration parameter, 682
  - transaction, 26
  - transaction ID
    - wraparound, 753
  - transaction isolation, 462
  - transaction isolation level, 463, 648
    - read committed, 463
    - repeatable read, 465
    - serializable, 466
    - setting, 1908
    - setting default, 647
  - transaction log (see [WAL](#))
  - transaction\_buffers configuration parameter, 595
  - transaction\_deferrable configuration parameter, 648
  - transaction\_isolation configuration parameter, 648
  - transaction\_read\_only configuration parameter, 648
  - transaction\_timeout configuration parameter, 649
  - transaction\_timestamp, 254
  - transform\_null\_equals configuration parameter, 659
  - transition tables, 1732
    - (see also [ephemeral named relation](#))
    - implementation in PLs, 1418
    - referencing from C trigger, 1236
  - translate, 217
  - TRANSLATE\_REGEX, 240
  - transparent huge pages, 592
  - trigger, 196, 1233
    - arguments for trigger functions, 1235
    - constraint trigger, 1733
    - for updating a derived tsvector column, 441
    - in C, 1236
    - in PL/pgSQL, 1316
    - in PL/Python, 1371
    - in PL/Tcl, 1344



- compared with rules, 1274
- trim, 212, 220
- trim\_array, 328
- trim\_scale, 206
- troubleshooting, 907
- true, 153
- trunc, 206, 274
- TRUNCATE, 1914
- trusted
  - PL/Perl, 1359
- tsm\_handler, 196
- tsquery (data type), 162
- tsquery\_phrase, 278, 439
- tsvector (data type), 161
- tsvector concatenation, 438
- tsvector\_to\_array, 278
- tsvector\_update\_trigger, 392
- tsvector\_update\_trigger\_column, 392
- ts\_debug, 279, 453
- ts\_delete, 277
- ts\_filter, 277
- ts\_headline, 277, 436
- ts\_lexize, 279, 456
- ts\_parse, 279, 455
- ts\_rank, 277, 435
- ts\_rank\_cd, 278, 435
- ts\_rewrite, 278, 439
- ts\_stat, 279, 442
- ts\_token\_type, 279, 456
- txid\_current, 367
- txid\_current\_if\_assigned, 367
- txid\_current\_snapshot, 367
- txid\_snapshot\_xip, 367
- txid\_snapshot\_xmax, 367
- txid\_snapshot\_xmin, 367
- txid\_status, 367
- txid\_visible\_in\_snapshot, 367
- type (see [data type](#))
- type cast, 38, 48
- typedef
  - in ECPG, 1021

## U

- UNESCAPE, 34, 36
- Unicode escape
  - in identifiers, 34
  - in string constants, 36
- Unicode normalization, 210, 211
- unicode\_assigned, 212
- unicode\_nul\_character\_replacement\_in\_jsonb configuration parameter, 659
- unicode\_version, 370
- UNION, 122
  - determination of result type, 407
- unique constraint, 63
- unistr, 217
- Unix domain socket, 919

- unix\_socket\_directories configuration parameter, 586
- unix\_socket\_group configuration parameter, 586
- unix\_socket\_permissions configuration parameter, 587
- unknown, 196
- UNLISTEN, 1916
- unnest, 329
  - for multirange, 334
  - for tsvector, 278
- unqualified name, 83
- updatable views, 1751
- UPDATE, 23, 105, 1917
  - RETURNING, 106
- update\_process\_title configuration parameter, 640
- updating, 105
- upgrading, 572
- UPLOAD\_MANIFEST, 2322
- upper, 212, 333, 334
  - and locales, 728
- upper\_inc, 333, 334
- upper\_inf, 333, 334
- UPSERT, 1826
- URI, 917
- usage\_tracking\_interval configuration parameter, 601
- user, 354, 710
  - current, 352
- user mapping, 101
- User name maps, 697
- user\_catalog\_table storage parameter, 1710
- USE\_BONJOUR, 2099
- USE\_BSD\_AUTH, 2099
- USE\_ICU, 2099
- USE\_LDAP, 2099
- USE\_LIBUNWIND, 2099
- USE\_LIBXML, 2100
- USE\_LIBXSLT, 2100
- USE\_LLVM, 2100
- USE\_LZ4, 2100
- USE\_OPENSSL, 2100
- USE\_PAM, 2100
- USE\_SYSTEMD, 2100
- USE\_ZSTD, 2100
- UUID, 163
  - generating, 163
- uuid-oss, 2764
- uuid\_extract\_timestamp, 279
- uuid\_extract\_version, 279
- uuid\_generate\_v1, 2764
- uuid\_generate\_v1mc, 2764
- uuid\_generate\_v3, 2764

## V

- vacuum, 750
- VACUUM, 1922
- vacuumdb, 2073
- vacuumlo, 2899



- `vacuum_buffer_usage_limit` configuration parameter, 594
- `vacuum_cost_delay` configuration parameter, 597
- `vacuum_cost_limit` configuration parameter, 598
- `vacuum_cost_page_dirty` configuration parameter, 598
- `vacuum_cost_page_hit` configuration parameter, 598
- `vacuum_cost_page_miss` configuration parameter, 598
- `vacuum_failsafe_age` configuration parameter, 650
- `vacuum_freeze_min_age` configuration parameter, 650
- `vacuum_freeze_table_age` configuration parameter, 649
- `vacuum_index_cleanup` storage parameter, 1708
- `vacuum_multixact_failsafe_age` configuration parameter, 650
- `vacuum_multixact_freeze_min_age` configuration parameter, 650
- `vacuum_multixact_freeze_table_age` configuration parameter, 650
- `vacuum_truncate` storage parameter, 1708
- `vacuum_update_datfrozenxid_only_when_needed` configuration parameter, 651
- value expression, 41
- VALUES, 124, 1927
  - determination of result type, 407
- varchar, 140
- variadic function, 1164
- variance, 339
  - population, 339
  - sample, 339
- `var_pop`, 339
- `var_samp`, 339
- version, 13, 369
  - compatibility, 572
- view, 25
  - implementation through rules, 1253
  - materialized, 1260
  - updating, 1266
- Visibility Map, 2467
- VM (see [Visibility Map](#))
- void, 196
- VOLATILE, 1172
- volatility
  - functions, 1172
- VPATH, 1232
- W**
- WAL, 845
- `wal_block_size` configuration parameter, 664
- `wal_buffers` configuration parameter, 605
- `wal_compression` configuration parameter, 605
- `wal_consistency_checking` configuration parameter, 668
- `wal_debug` configuration parameter, 668
- `wal_decode_buffer_size` configuration parameter, 609
- `wal_init_zero` configuration parameter, 605
- `wal_keep_size` configuration parameter, 613
- `wal_level` configuration parameter, 602
- `wal_log_hints` configuration parameter, 605
- `wal_receiver_create_temp_slot` configuration parameter, 617
- `wal_receiver_status_interval` configuration parameter, 617
- `wal_receiver_timeout` configuration parameter, 617
- `wal_recycle` configuration parameter, 605
- `wal_retrieve_retry_interval` configuration parameter, 617
- `wal_segment_size` configuration parameter, 664
- `wal_sender_check_crc` configuration parameter, 606
- `wal_sender_panic_on_crc_error` configuration parameter, 606
- `wal_sender_timeout` configuration parameter, 613
- `wal_skip_threshold` configuration parameter, 606
- `wal_summary_keep_time` configuration parameter, 612
- `wal_sync_method` configuration parameter, 604
- `wal_writer_delay` configuration parameter, 605
- `wal_writer_flush_after` configuration parameter, 606
- warm standby, 770
- `websearch_to_tsquery`, 276
- WHERE, 115
- where to log, 627
- WHILE
  - in PL/pgSQL, 1301
- width, 269
- `width_bucket`, 207
- window function, 27
  - built-in, 341
  - invocation, 46
  - order of execution, 120
- `wipe_file_on_delete` configuration parameter, 660
- `wipe_heaptuple_on_delete` configuration parameter, 660
- `wipe_memctx_on_free` configuration parameter, 660
- `wipe_mem_on_free` configuration parameter, 660
- `wipe_xlog_on_free` configuration parameter, 660
- WITH
  - in SELECT, 125, 1878
- WITH CHECK OPTION, 1749
- WITHIN GROUP, 44
- witness server, 770
- `work_mem` configuration parameter, 593
- wraparound
  - of multixact IDs, 755
  - of transaction IDs, 753
- `write_page_cost` configuration parameter, 623

**X**

- xid, 193
- xid8, 193
- xmax, 69
- xmin, 68
- XML, 163
- XML export, 289
- XML Functions, 280
- XML option, 164, 651
- xmlagg, 283, 337
- xmlbinary configuration parameter, 651
- xmlcomment, 280
- xmlconcat, 280
- xmlelement, 281
- XMLEXISTS, 284
- xmlforest, 282
- xmloption configuration parameter, 651
- xmlparse, 164
- xmlpi, 283
- xmlroot, 283
- xmlserialize, 164
- xmltable, 286
- xmltext, 280
- xml\_is\_well\_formed, 284
- xml\_is\_well\_formed\_content, 284
- xml\_is\_well\_formed\_document, 284
- xml\_parse\_huge configuration parameter, 592
- XPath, 285
- xpath\_exists, 286
- XQuery regular expressions, 240

**Z**

- zero\_damaged\_pages configuration parameter, 668